



THE UNIVERSITY *of* EDINBURGH

## Edinburgh Research Explorer

### Computer Aided Translation: Advances and Challenges

**Citation for published version:**

Koehn, P 2015, 'Computer Aided Translation: Advances and Challenges', MT Summit XV, Miami, FL, United States, 30/10/15 - 3/11/15. <[https://amtaweb.org/wp-content/uploads/2016/10/Computer-Aided\\_Tutorial\\_Koehn\\_wide-cover.pdf](https://amtaweb.org/wp-content/uploads/2016/10/Computer-Aided_Tutorial_Koehn_wide-cover.pdf)>

**Link:**

[Link to publication record in Edinburgh Research Explorer](#)

**Document Version:**

Publisher's PDF, also known as Version of record

**General rights**

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

**Take down policy**

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact [openaccess@ed.ac.uk](mailto:openaccess@ed.ac.uk) providing details, and we will remove access to the work immediately and investigate your claim.



---

# Computer Aided Translation Advances and Challenges

Philipp Koehn

30 October 2015



- A practical introduction: the CASMACAT workbench
  - Postediting
  - Types of assistance
  - Logging, eye tracking and user studies
  - Implementation details of the CASMACAT workbench
- :

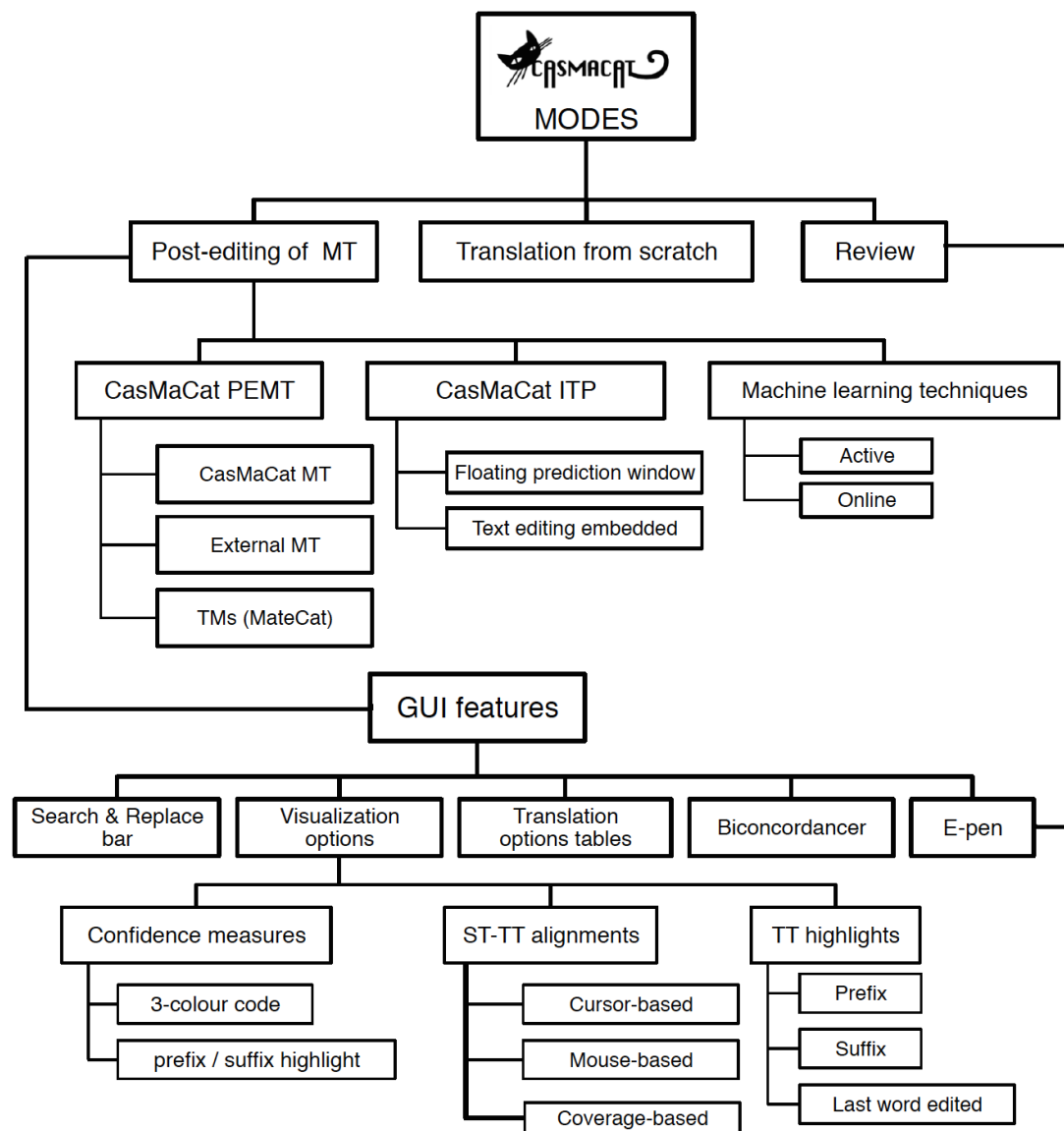
# part I

## CASMACAT workbench

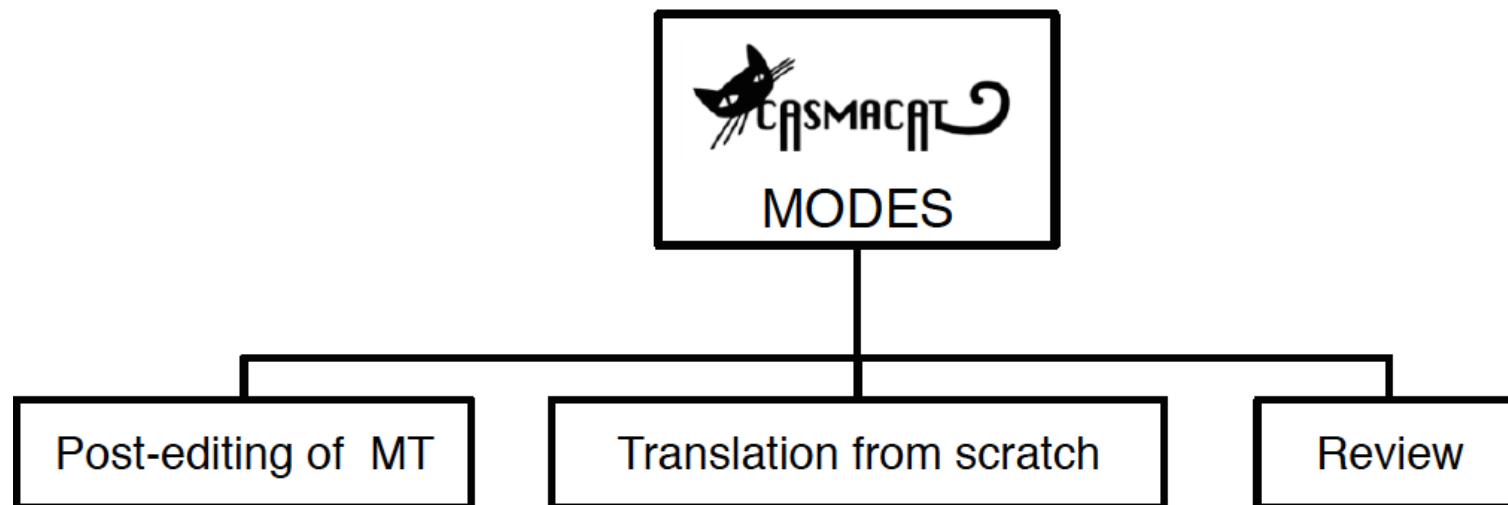
# CASMACAT workbench

- Cognitive studies of translators leading to insights into interface design
  - better understanding of translator needs
- Workbench with novel types of assistance to human translators
  - interactive translation prediction
  - interactive editing and reviewing
  - adaptive translation models
  - better tools for translators
- Demonstration of effectiveness in field tests with professional translators
  - increased translator productivity

# Architecture



# Core Modes

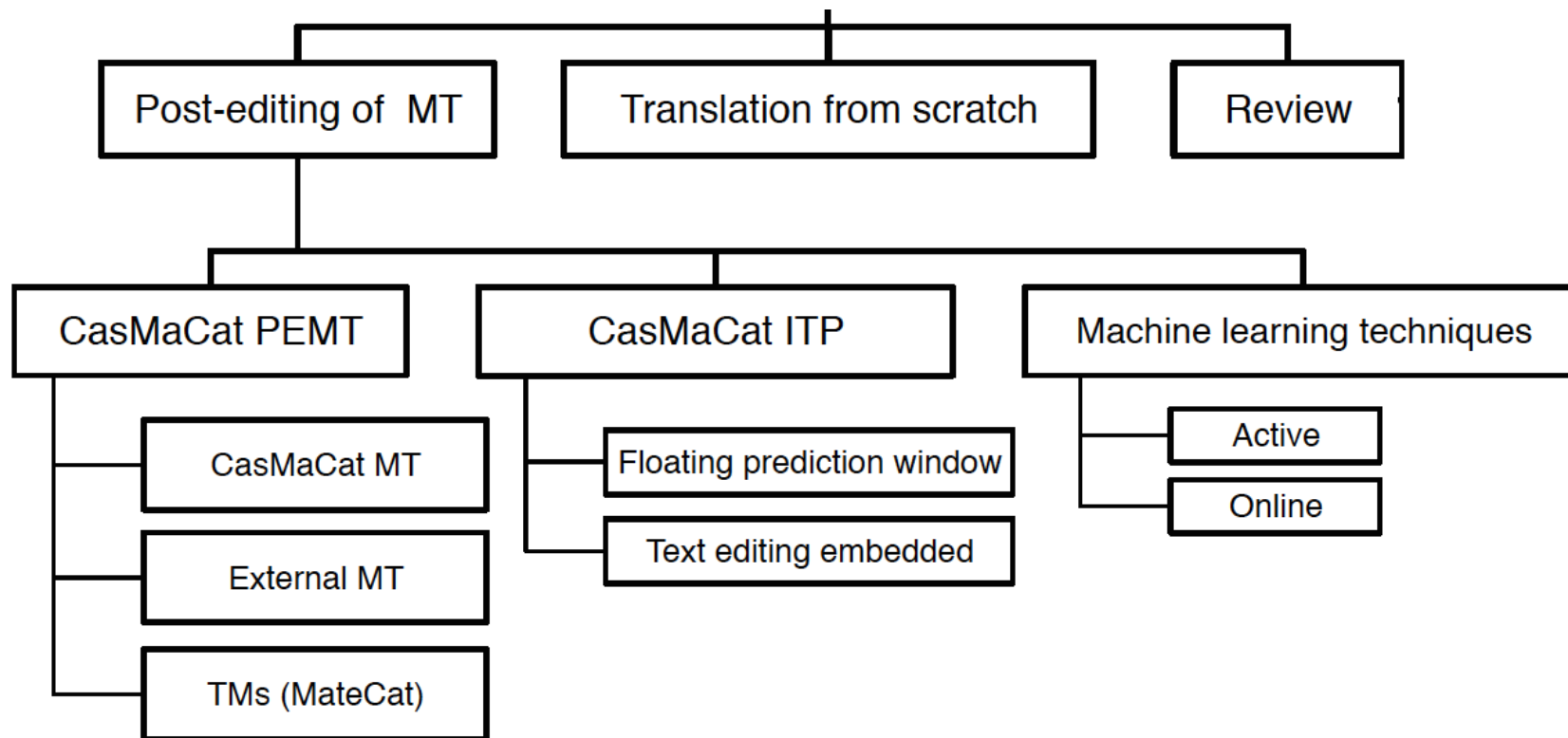


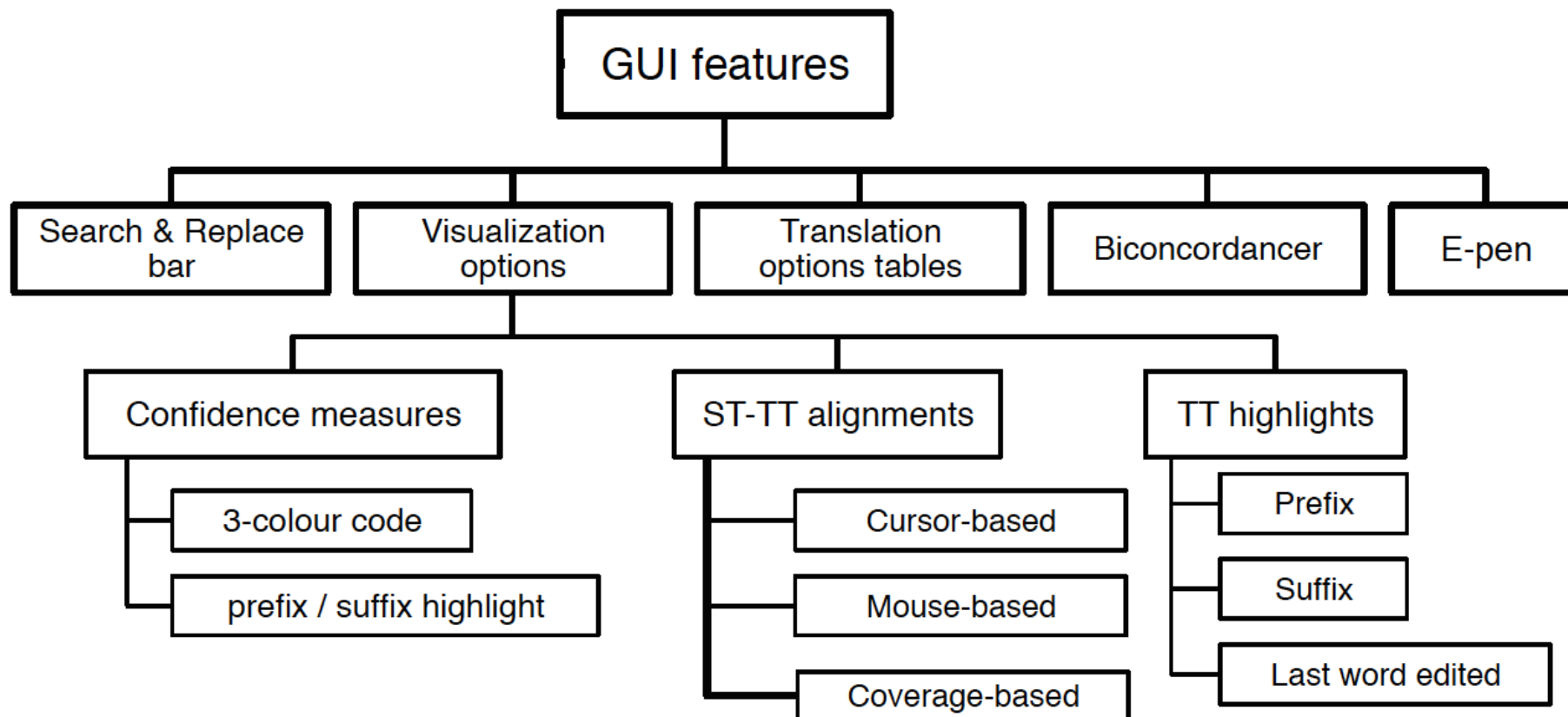


# Postediting Modes



7





# Postediting Interface



6 Le Pakistan a donc été récompensé par l'assistance et les armes des États-Unis. > As a result, Pakistan was rewarded with American financial assistance and arms.

7 Pour mieux redistribuer ses cartes, Moucharraaf a envoyé l'armée pakistanaise dans les zones ethniques qui longent l'Afghanistan, pour la première fois depuis l'indépendance du Pakistan. > In furtherance of his re-alignment, Musharraf sent the Pakistani army into the tribal areas bordering Afghanistan for the first time since Pakistan's independence.

8 Les opérations contre les forces des Talibans et d'Al-Qaeda ont obtenu des résultats mitigés. >

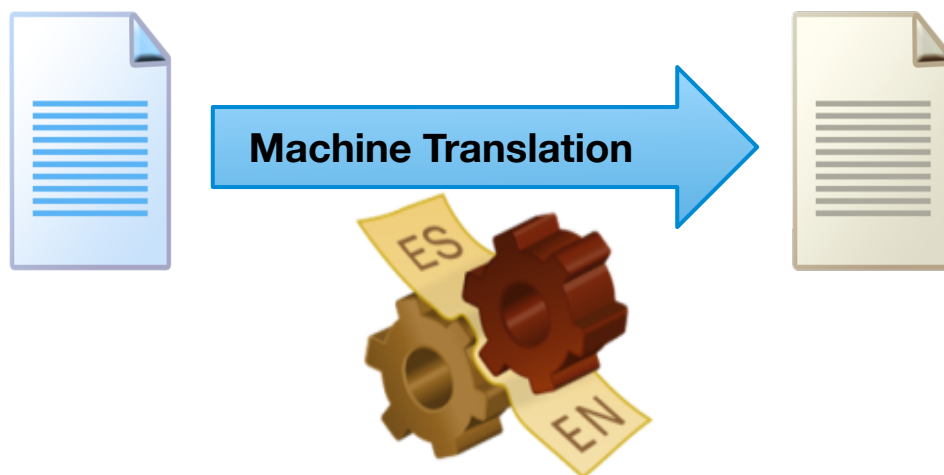
visualization >>

ITP T→ DRAFT **TRANSLATED**

- Source on left, translation on right
- Context above and below

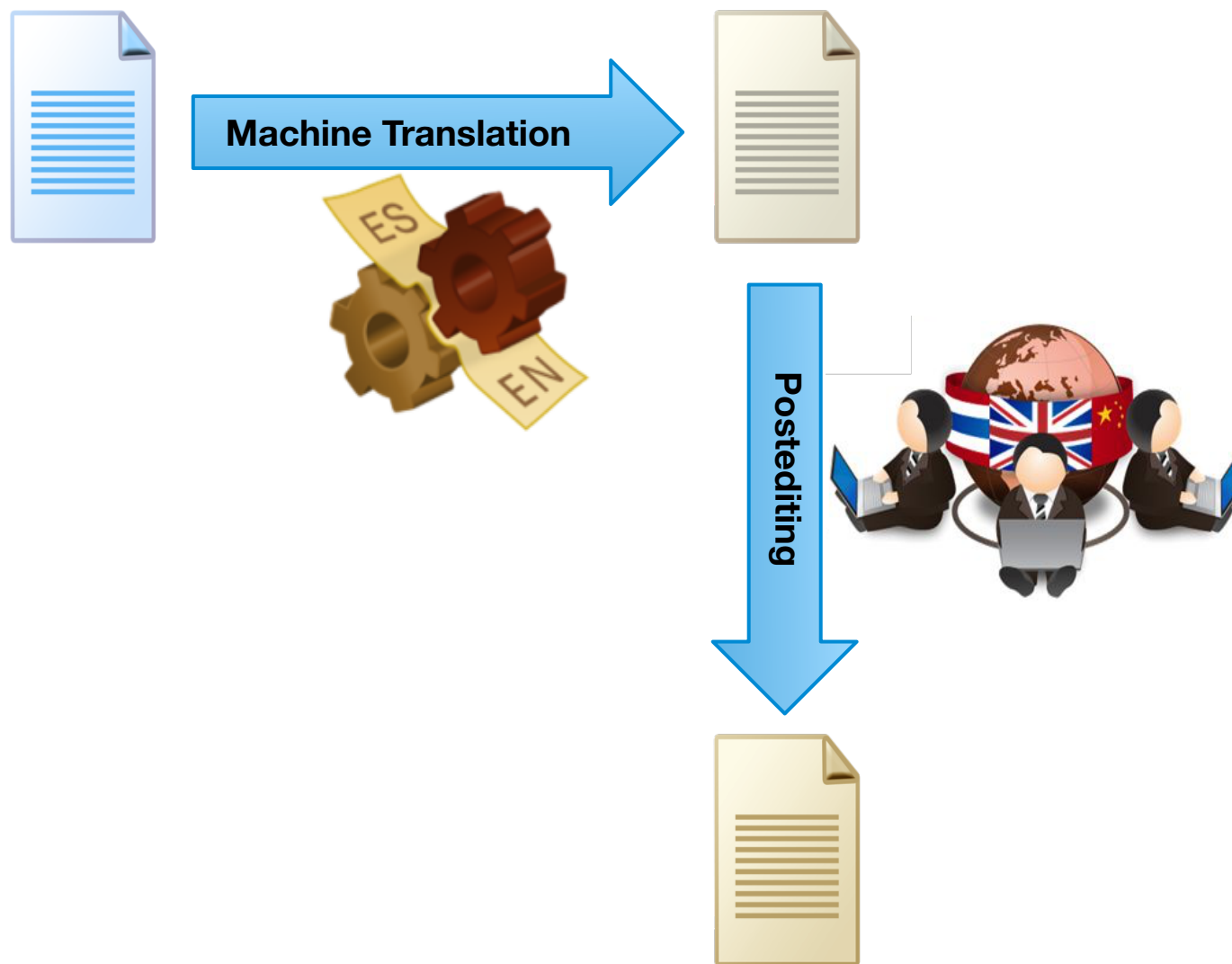
# Incremental Updating

10



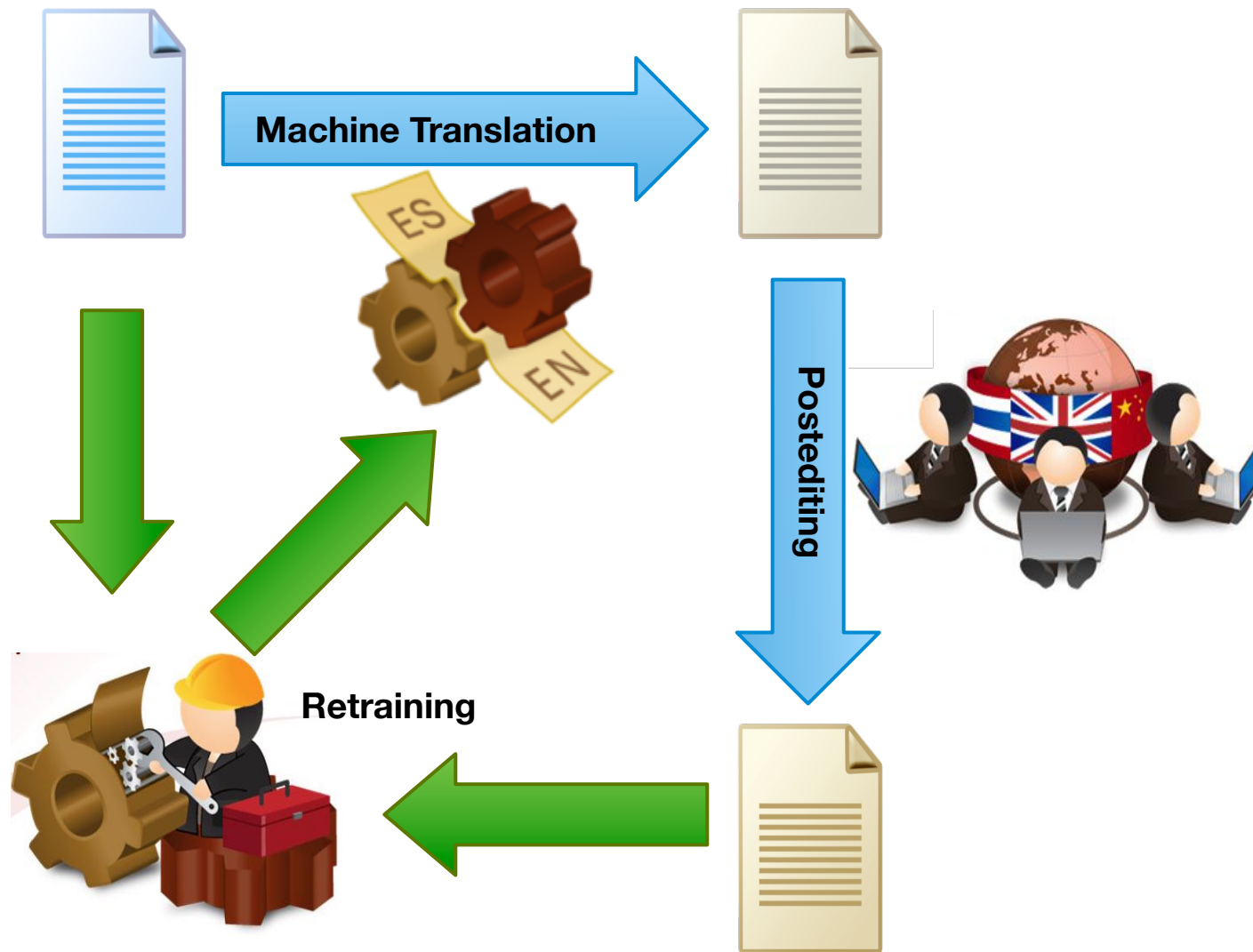
# Incremental Updating

11



# Incremental Updating

12



# Word Alignment

13



6 Le Pakistan a donc été récompensé par l'assistance et les armes des États-Unis. As a result, Pakistan was rewarded with American financial assistance and arms.

7 **visualization >>** ☒ displayMouseAlign ☒ displayCaretAlign ☐ displayShadeOffTranslatedSource ☐ displayConfidences ☐ highlightValidated ☐ highlightPrefix ☐ highlightLastValidated ☐ limitSuffixLength

Pour mieux redistribuer ses cartes, Moucharraaf a envoyé l'armée pakistanaise dans les **zones ethniques** qui **longent** l'**Afghanistan**, pour la première fois depuis l'indépendance du Pakistan.

In furtherance of his re-alignment, Musharraf sent the Pakistani army into the **tribal** areas bordering Afghanistan for the first time since Pakistan's independence.

ITP T→ DRAFT **TRANSLATED**

8 Les opérations contre les forces des Talibans et d'Al-Qaeda ont obtenu des résultats mitigés.

- Caret alignment (green)
- Mouse alignment (yellow)



The screenshot shows a machine translation interface. On the left, the source text is "And on that the signs are mixed." On the right, the target text is "Y en que los indicios son desiguales." The target text has color-coded words: "en" is orange, "los" is orange, "indicios" is orange, "son" is red, and "desiguales" is red. Below the source text, there is a tab labeled "Translation matches" and the same source text. Below the target text, there is a status bar with the source "ITP", the date "Fri Apr 12 2013 18:03:17 GMT+0200 (CEST)", and a score of "42". Above the target text, there are buttons for "ITP", "T→", "DRAFT", and "TRANSLATED".


- Sentence-level confidence measures  
→ estimate usefulness of machine translation output
- Word-level confidence measures  
→ point posteditor to words that need to be changed



# Interactive Translation Prediction

15





Re-calibrateDownload edf-fileDOWNLOAD PROJECTHELP

Document list > Jobs List > fiction.xliff...fiction.xliff  
Shortcuts(29) > en-GB > es-ES



10314

Forget it. It's too risky. I'm through doing that shit.

visualization >>

Olvidarlo. Es demasiado

arriesgado. Estoy haciendo

ITPT→DRAFTTRANSLATED

10315

You always say that. The same thing every time.

10316

"I'm through, never again, too dangerous."

# Bilingual Concordancer

16



TIP ≡ T→ DRAFT TRANSLATED

abandonner

✕

**abandon**

ances des Etats-Unis à	<b>abandonner</b>	Musharraf -- et les co		merican reluctance to	<b>abandon</b>	Musharraf -- together
uridique, il a décidé d'	<b>abandonner</b>	la constitutionnalité, c		af has now decided to	<b>abandon</b>	constitutionality, remc
implement menacé d'	<b>abandonner</b>	ses accords commerci		simply threatened to	<b>abandon</b>	or never to conclude t

**give up**

erait donc contraint d'	<b>abandonner</b>	le droit de créer son p		e would be required to	<b>give up</b>	the right to develop it
n' était pas disposé à	<b>abandonner</b>	ses fonctions militaire		arraf was not ready to	<b>give up</b>	his military post, but a

**to**

t ne veulent donc pas	<b>abandonner</b>	leurs prérogatives dar		olicy and do not want	<b>to</b>	delegate this prerogat
-----------------------	-------------------	------------------------	--	-----------------------	-----------	------------------------

**to abandon**

es tout en refusant d'	<b>abandonner</b>	son arsenal nucléaire		drawal while refusing	<b>to abandon</b>	its nuclear weapons a
------------------------	-------------------	-----------------------	--	-----------------------	-------------------	-----------------------

# Translation Option Array

17



hikers are severely injured, and ten people are missing.  
 after Mount Ontake (御嶽山, Ontake-san), a popular climbing  
 spot in central Japan, **erupted** for the first time in five years.

Kletterer sind schwer verletzt, und zehn Menschen werden  
 vermisst, nachdem Mount Ontake (御嶽山, Ontake-san), ein  
 beliebter Kletterplatz im zentralen Japan, |

**ausbruch, zum ersten**

ITP ≡ T→ DRAFT **TRANSLATED**

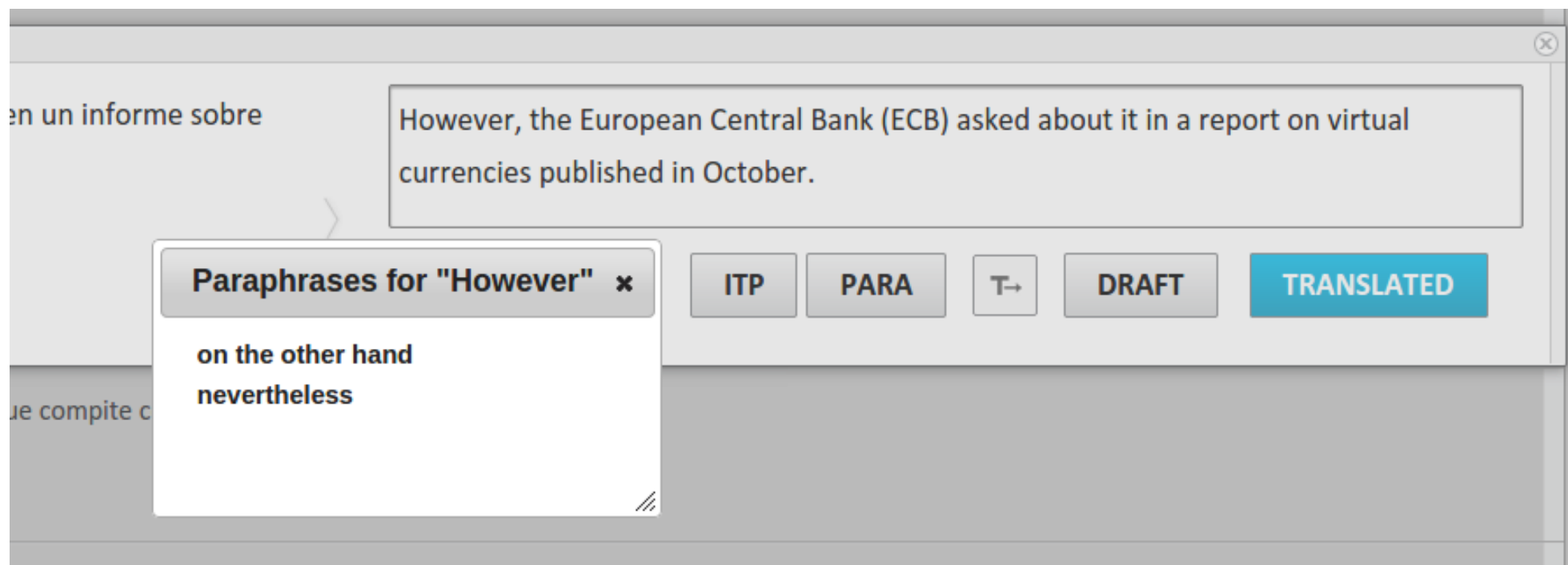
Translation Options

ke	-	san	) ,	a	popular	climbing	spot	in central	Japan ,	<b>erupted</b>	for the first time in five years .
ke	-	san	),	ein	beliebtes	Klettern	vor Ort	in Mittel-	Japan,	ausbrach	zum ersten Mal in fünf Jahren .
	und	San	) ,	ein	populär	Bergsteigen	vor	zentrale	Japan ,	ausbrach,	zum ersten Mal in fünf Jahre.
	/		), die		beliebt	Aufstieg	Fleck	zentralen	Japans,	platzte	zum ersten Mal fünf Jahre
	der		)	eine	beliebte	abhalten,	ein, in	zentraler	Japan	Ausbruch	in fünf Jahren
	bis		), in	populär		Erklimmen	Vor - Ort @-@	zentral	Japans .	ausgebrochen	zum ersten Mal in der von fünf Jahren.
	von		), .	populär ist,		beim Besteigen	in	mittel-	in Japan -	ausgebrochen ist	zum ersten Mal seit fünf Jahren sind.

- Visual aid: non-intrusive provision of cues to the translator
- Clickable: click on target phrase → added to edit area
- Automatic orientation
  - most relevant is next word to be translated
  - automatic centering on next word

# Paraphrasing

18

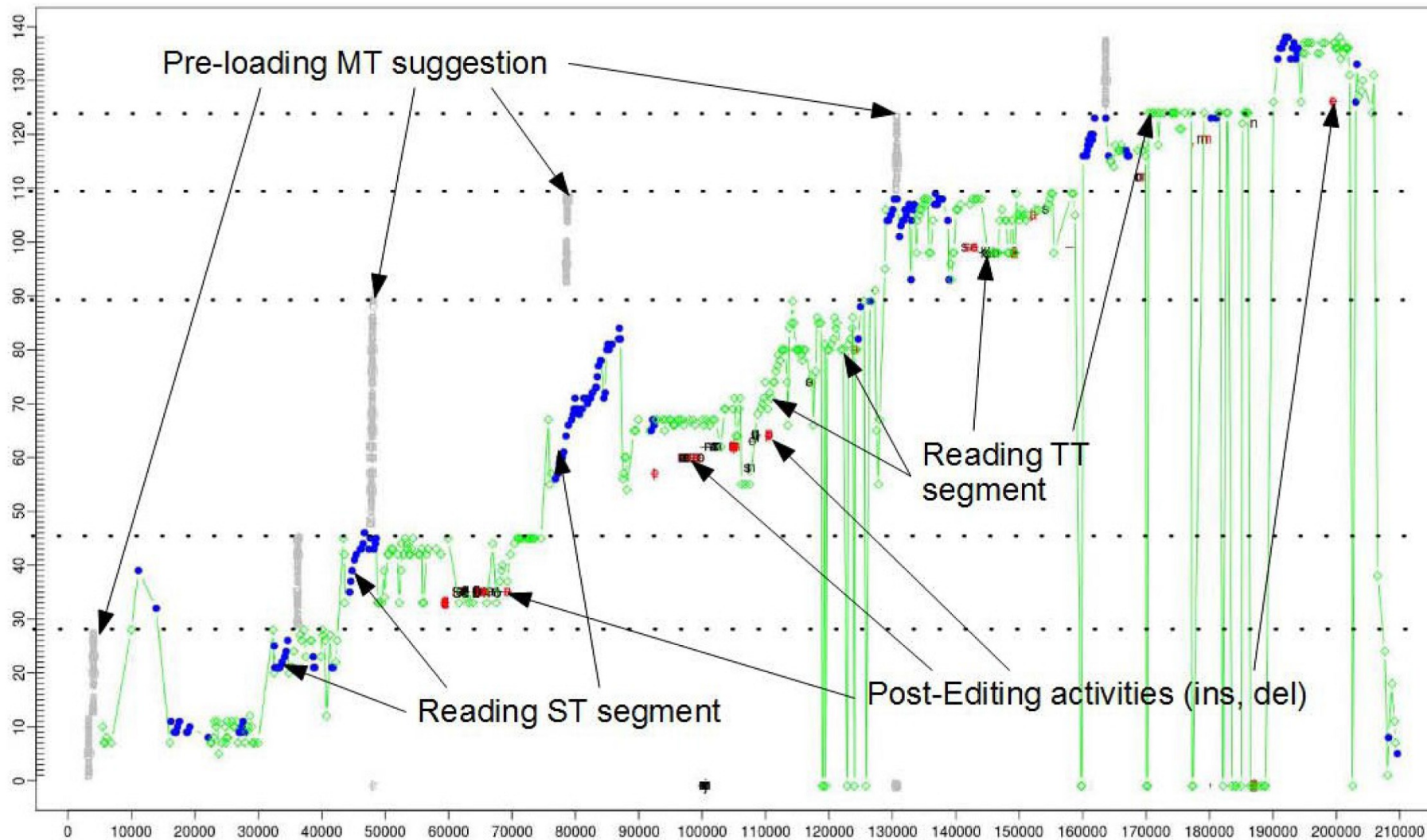


# How do we Know it Works?

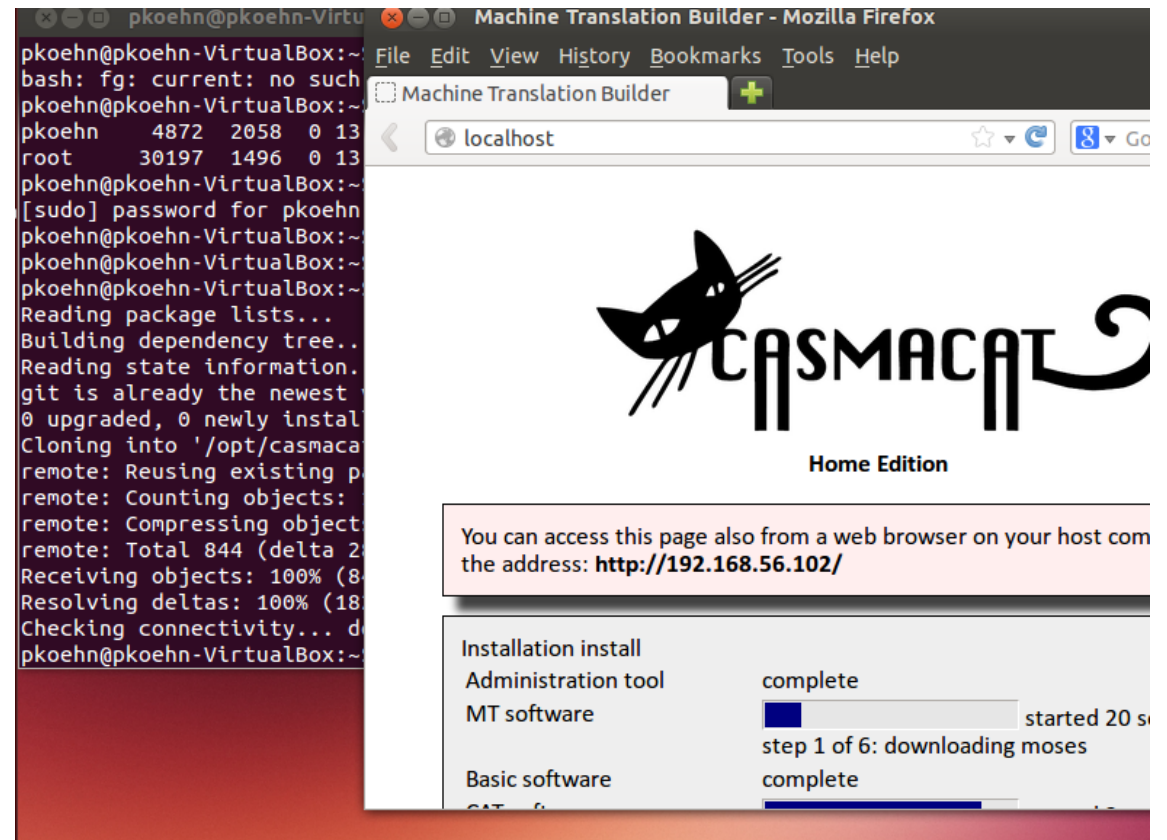
- Intrinsic Measures
  - word level confidence: user does not change words generated with certainty
  - interactive prediction: user accepts suggestions
- User Studies
  - professional translators faster with post-editing
  - ... but like interactive translation prediction better
- Cognitive studies with eye tracking
  - where is the translator looking at?
  - what causes the translator to be slow?

# Logging and Eye Tracking

20



- Running CSMACAT on your desktop or laptop
- Installation
  - Installation software to run virtual machines (e.g., Virtualbox)
  - installation of Linux distribution (e.g., Ubuntu)
  - installation script sets up all the required software and dependencies



# Administration through Web Browser



## Administration

### Translate

- [Translate new document](#)
- [List documents](#)

### Engines

- [Manage engines](#)
- [Upload engine](#)
- [Build new prototype](#)

### Settings

- [Reset CAT and MT server](#)
- [CAT Settings](#)
- [Update Software](#)

**Deployed:** fr-en-upload-1  
**Memory:** 1.2 GB used, 6.6 GB free  
**Disk:** 12.9 GB used, 10.2 GB free  
**Uptime:** 22:24  
**Load:** 0.01, 0.05, 0.08  
Monday, 06 October 2014, 21:22:41





# Training MT Engines

23



- Train MT engine on own or public data

### Build New Prototype

Input language

Output language

Add corpus  No file chosen

Name	Segments	Publisher	
<a href="#">European Central Bank</a>	102,980	OPUS	<a href="#">upload</a>
<a href="#">European Medicines Agency</a>	372,824	OPUS	<a href="#">upload</a>
<a href="#">EU Bookshop</a>	3,618,897	OPUS	<a href="#">upload</a>
<a href="#">European Constitution</a>	6,667	OPUS	<a href="#">upload</a>
<a href="#">European Parliament</a>	1,260,689	OPUS	<a href="#">upload</a>
<a href="#">KDE4</a>	126,141	OPUS	uploaded
<a href="#">KDE4 (el-en_GB)</a>	125,537	OPUS	<a href="#">upload</a>
<a href="#">Open Subtitles</a>	220,445	OPUS	<a href="#">upload</a>
<a href="#">Open Subtitles 2011</a>	10,693,456	OPUS	<a href="#">upload</a>
<a href="#">Open Subtitles 2012</a>	12,984,773	OPUS	<a href="#">upload</a>
<a href="#">Open Subtitles 2013</a>	14,626,890	OPUS	<a href="#">upload</a>
<a href="#">South-East European Times</a>	165,532	OPUS	<a href="#">upload</a>
<a href="#">South-East European Times v2</a>	224,808	OPUS	<a href="#">upload</a>
<a href="#">SPC</a>	7,035	OPUS	<a href="#">upload</a>
<a href="#">Tatoeba</a>	2,469	OPUS	<a href="#">upload</a>
<a href="#">DGT-Translation Memory</a>	3,016,402	JRC	<a href="#">upload</a>

**Corpora**

Use	ID	Name	Segments	Uploaded
<input checked="" type="checkbox"/> all	1	KDE4	126141	21:39:27

**Re-Use** Previous setting

**Tuning set**  ☐ all ☒ select

**Evaluation set**  ☐ all ☒ select

**Name**

# Managing MT Engines

24



- MT engines can be
  - switched out
  - downloaded
  - uploaded
  - shared

## Manage Engines

### English-French

Available Engines

#	Name	Size	Build date	Action
2	NC+TED	2.3G	27 Mar 14	<a href="#">deploy</a> <a href="#">delete</a> <a href="#">download</a>

Prototypes ([Inspect Details in Prototype Factory](#))

#	Name	Status	Build date	Action
2	NC+TED	done	Fri 20:34	<a href="#">delete</a>
1	NC	done	Fri 20:34	<a href="#">create engine</a> <a href="#">delete</a>

### English-Spanish

Available Engines

#	Name	Size	Build date	Action
2	NC+TED	2.3G	27 Mar 14	<a href="#">deploy</a> <a href="#">delete</a> <a href="#">download</a>

Prototypes ([Inspect Details in Prototype Factory](#))

#	Name	Status	Build date	Action
3	NC+TED+EP	stopped	Fri 20:34	<a href="#">resume</a> <a href="#">delete</a>
2	NC+TED	done	Fri 20:34	<a href="#">delete</a>
1	NC	done	Fri 20:34	<a href="#">create engine</a> <a href="#">delete</a>

### French-English

Available Engines

#	Name	Size	Build date	Action
x1	Toy	85M	27 Mar 14	deployed <a href="#">download</a>
2	NC+TED	2.3G	27 Mar 14	<a href="#">deploy</a> <a href="#">delete</a> <a href="#">download</a>

Prototypes ([Inspect Details in Prototype Factory](#))

#	Name	Status	Build date	Action
2	NC+TED	done	Fri 20:34	<a href="#">delete</a>
1	NC	done	Fri 20:34	<a href="#">create engine</a> <a href="#">delete</a>

### Spanish-English

Available Engines

#	Name	Size	Build date	Action
2	NC+TED	2.3G	27 Mar 14	<a href="#">deploy</a> <a href="#">delete</a> <a href="#">download</a>

# CAT Settings

- With own MT engine, all CASKMACAT modes are available

**CAT Settings**

Interactive Translation Prediction	<input checked="" type="checkbox"/>
Search and Replace	<input checked="" type="checkbox"/>
Bilingual Concordancer	<input type="checkbox"/>
Hide Contributions	<input checked="" type="checkbox"/>
Floating Predictions	<input checked="" type="checkbox"/>
Translation Options	<input type="checkbox"/>
Allow Change of Visualization Options	<input checked="" type="checkbox"/>
Restrict ITP to Draft Stage	<input type="checkbox"/>

update



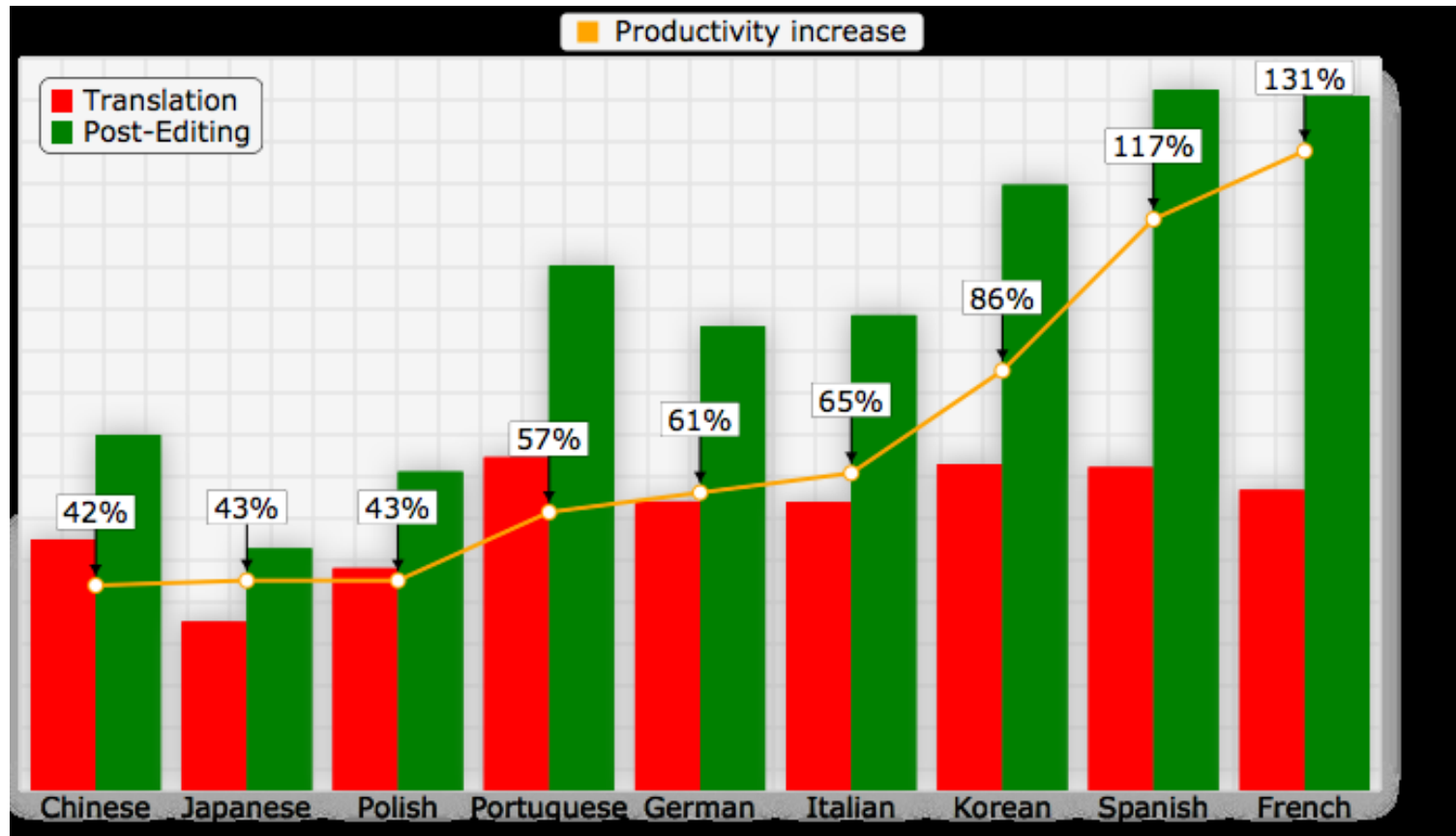
# part II

## cat methods



# post-editing

# Productivity Improvements



(source: Autodesk)

# MT Quality and Productivity

- What is the relationship between MT Quality and Postediting Speed
- One study (English–German, news translation, non-professionals)

System	Speed		Metric	
	sec./wrd.	wrds./hr.	BLEU	MANUAL
ONLINE-B	5.46	659	20.7	0.637
UEDIN-SYNTAX	5.38	669	19.4	0.614
UEDIN-PHRASE	5.45	661	20.1	0.571
UU	6.35	567	16.1	0.361

# Translator Variability

- Translator differ in
  - ability to translate
  - motivation to fix minor translation
- High variance in translation time  
(again: non-professionals)

Post-editor	Speed	
	sec./wrd.	wrds./hr.
1	3.03	1,188
2	4.78	753
3	9.79	368
4	5.05	713



# MT Quality and Postediting Effort

- Postediting effort = number of words changed
- Evaluation metric at IWSLT 2014
  - TER = automatic metric, comparison against a reference translation
  - HTER = postediting metric, actual words changed

## English–German

Ranking	HTER	TER
EU-BRIDGE	19.2	54.6
UEDIN	19.9	56.3
KIT	20.9	54.9
NTT-NAIST	21.3	54.7
KLE	28.8	59.7

## English–French

Ranking	HTER	TER
EU-BRIDGE	16.5	42.6
RWTH	16.6	41.8
KIT	17.6	42.3
UEDIN	17.2	43.3
MITLL-AFRL	18.7	43.5
FBK	22.3	44.3
MIRACL	32.9	52.2

- Professional translators

## English–German

Posteditor	HTER	TER
PE 1	32.2	56.1
PE 2	19.7	56.3
PE 3	40.9	56.2
PE 4	27.6	55.9
PE 5	25.0	55.6

## English–French

Posteditor	HTER	TER
PE 1	35.0	42.6
PE 2	17.5	42.8
PE 3	23.7	43.0
PE 4	39.7	42.3
PE 5	19.7	42.9

- Also very high variability

- Goal of MT quality metrics not clear
  - understandability: do you get the meaning?
  - post-editing effort: how much effort to change?
- Example: dropping of the word "not"
  - understandability: big mistake
  - post-editing effort: quick add of just one word
- Not clear, what tradition manual metrics prefer (adequacy, fluency)
- Not clear, what BLEU score etc. prefer



# word alignment

# Word Alignment

35



6 Le Pakistan a donc été récompensé par l'assistance et les armes des États-Unis. As a result, Pakistan was rewarded with American financial assistance and arms.

7 **visualization >>** ☒ displayMouseAlign ☒ displayCaretAlign ☐ displayShadeOffTranslatedSource ☐ displayConfidences ☐ highlightValidated ☐ highlightPrefix ☐ highlightLastValidated ☐ limitSuffixLength

Pour mieux redistribuer ses cartes, Moucharraaf a envoyé l'armée pakistanaise dans les **zones ethniques** qui **longent** l'**Afghanistan**, pour la première fois depuis l'indépendance du Pakistan.

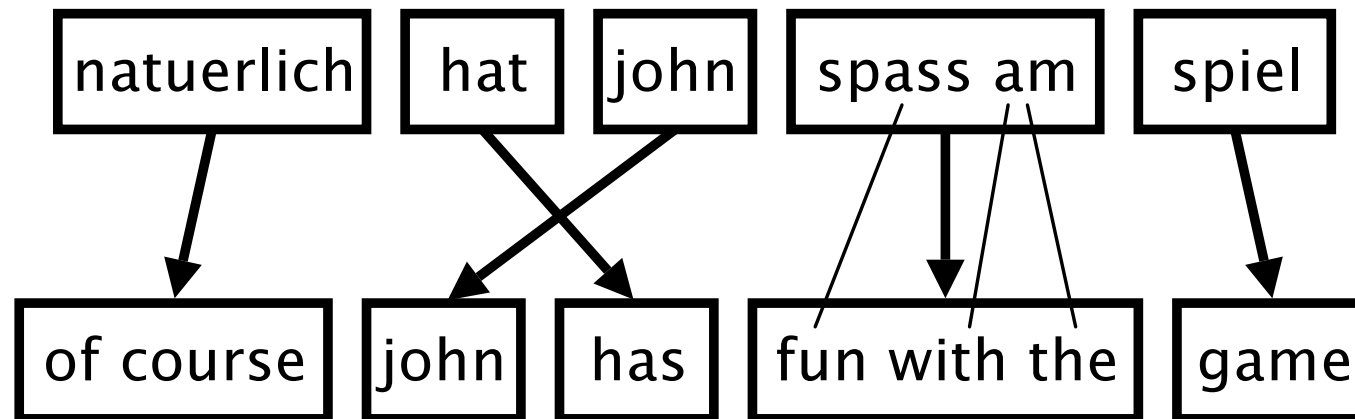
In furtherance of his re-alignment, Musharraf sent the Pakistani army into the **tribal** areas bordering Afghanistan for the first time since Pakistan's independence.

ITP T→ DRAFT **TRANSLATED**

8 Les opérations contre les forces des Talibans et d'Al-Qaeda ont obtenu des résultats mitigés.

- Caret alignment (green)
- Mouse alignment (yellow)

# Word Alignment from MT



- Machine translation output is constructed by phrase mappings
- Each phrase mapping has internal word alignment

⇒ This can be used to visualize word alignments

- But: word alignment points become invalid after user edits



- During machine translation training, standard component is word alignment
- Standard tools
  - old workhorse: GIZA++
  - currently popular tool: fast-align
- These tools have been adapted to align new sentence pairs

# Mouse Over Alignment

Pour mieux redistribuer ses cartes, Moucharraaf a envoyé l'armée pakistanaise dans les **zones ethniques** qui **longent** l'Afghanistan, pour la première fois depuis l'indépendance du Pakistan.



In furtherance of his re-alignment, Musharraaf sent the Pakistani army into the **tribal** areas bordering Afghanistan for the first time since Pakistan's independence.

- Highlight the source word aligned to the word at the current **mouse** position



# Caret Alignment

Pour mieux redistribuer ses cartes, Moucharraf a envoyé l'armée pakistanaise dans les zones ethniques qui longent l'Afghanistan, pour la première fois depuis l'indépendance du Pakistan.

In furtherance of his re-alignment, Musharraf sent the Pakistani army into the tribal areas bordering Afghanistan for the first time since Pakistan's independence.

- Highlight the source word aligned to the word at the current **caret** position

# Shade Off Translated

L'intervention israélienne dans la bande de Gaza et les bombardements américains en Irak pour lutter contre les djihadistes de l'État islamique en Irak et au Levant ont également ajouté de la nervosité sur les marchés.

Israeli intervention in the Gaza Strip and the

American bombing in

- Use in interactive prediction mode
- Shade off words that are already translated
- Highlight words aligned to first predicted translation word



# confidence measures

- Machine translation engine indicates where it is likely wrong
- Different Levels of granularity
  - document-level (SDL's "TrustScore")
  - sentence-level
  - word-level■
- What are we predicting?
  - how useful is the translation — on a scale of (say) 1–5
  - indication if post-editing is worthwhile
  - estimation of post-editing effort
  - pin-pointing errors

- Translators are used to "Fuzzy Match Score"
  - used in translation memory systems
  - roughly: ratio of words that are the same between input and TM source
  - if less than 70%, then not useful for post-editing
- We would like to have a similar score for machine translation■
- Even better
  - estimation of post-editing time
  - estimation of from-scratch translation time
  - can also be used for pricing
- Very active research area

# Quality Estimation Shared Task

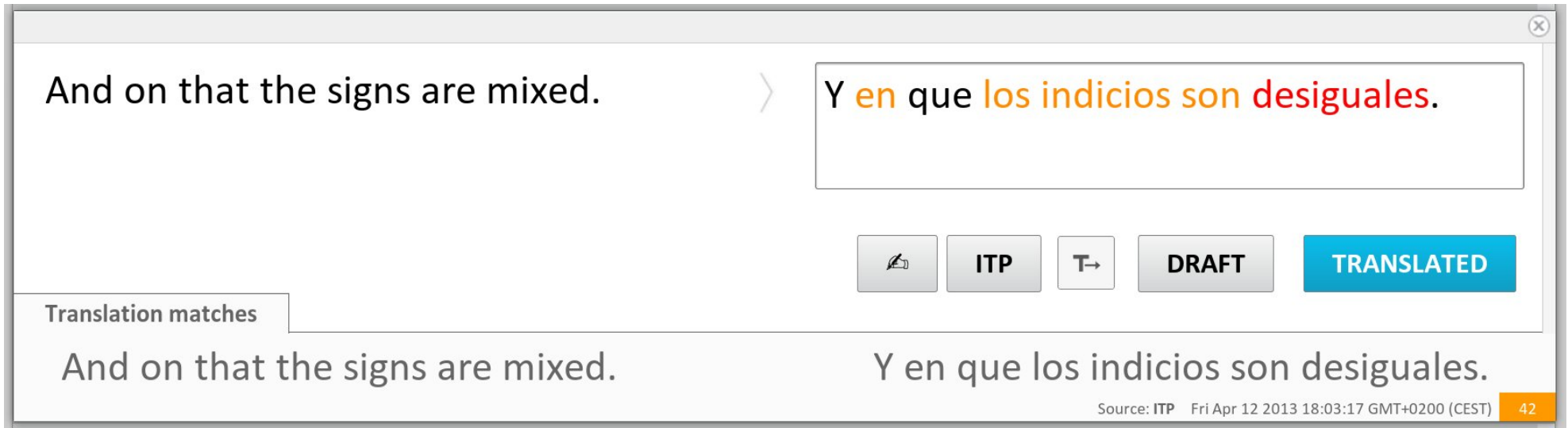
- Shared task organized at WMT since 2012
- Given
  - source sentence
  - machine translation
- Predict
  - human judgement of usefulness for post-editing (2012, 2014)
  - HTER score on post-edited sentences (2013, 2014, 2015)
  - post-editing time (2013, 2014)
- Also task for word-level quality estimation (2014, 2015) and document-level quality estimation (2015)

- Open source tool for quality estimation
- Source sentence features
  - number of tokens
  - language model (LM) probability
  - 1–3-grams observed in training corpus
  - average number of translations per word
- Similar target sentence features
- Alignment features
  - difference in number of tokens and characters
  - ratio of numbers, punctuation, nouns, verbs, named entities
  - syntactic similarity (POS tags, constituents, dependency relationships)
- Scores and properties of the machine translation derivation
- Uses Python's SCIKIT-LEARN implementation of SVM regression



# word level confidence





The screenshot displays a CAT interface with a main workspace and a sidebar. The main workspace shows a source text "And on that the signs are mixed." on the left and a target text "Y en que los indicios son desiguales." on the right. The target text is highlighted in orange. Below the target text, there are buttons for "ITP", "T→", "DRAFT", and "TRANSLATED". The sidebar on the left is titled "Translation matches" and shows a list of matches, including the source text "And on that the signs are mixed." and the target text "Y en que los indicios son desiguales.".

And on that the signs are mixed. > Y en que los indicios son desiguales.

Translation matches

And on that the signs are mixed. Y en que los indicios son desiguales.

Source: ITP Fri Apr 12 2013 18:03:17 GMT+0200 (CEST) 42

- Highlight words less likely to be correct

- Simple methods quite effective
  - IBM Model 1 scores
  - posterior probability of the MT model
- Machine learning approach
  - similar features as for sentence-level quality estimation

- Machine translation output

*Quick brown fox jumps on the dog lazy.*

- Post-editing

*The quick brown fox jumps over the lazy dog.*

- Annotation

<i>Fast</i>	<i>brown</i>	<i>fox</i>	<i>jumps</i>	<i>on</i>	<i>the</i>	<i>dog</i>	<i>lazy</i>	<i>.</i>
bad	good	good	good	bad	good	good	good	good

- Problems: dropped words? reordering?

- Evaluated in user study
- Feedback
  - could be useful feature
  - but accuracy not high enough
- To be truly useful, accuracy has to be very high
- Current methods cannot deliver this



# automatic reviewing

- Can we identify errors in human translations?
  - missing / added information
  - inconsistent use of terminology

## Input Sentence

Er hat seit Monaten geplant, im Oktober einen Vortrag in Miami zu halten.

## Human Translation

Moreover, he planned for months to give a lecture in Miami.

- Intuition
  - reviewing more efficient with pen and paper
  - e-pen enables this work process in digital environment
- Work carried out
  - fronted modified for larger drawing area
  - backend support for hand-written text recognition (HTR)
  - development of methods for HTR
- Field trial carried out → corpus of reviewing edits

# Analysis of Reviewer Edits

- 171 insertions — vast majority function words
- 152 deletions — about half substantial content
- 621 replacements — of which:
  - 75 changes to punctuation only
  - 28 change to lowercase / uppercase
  - 29 cases that are mostly deletions
  - 8 cases that are mostly insertions
  - 289 morphological/spelling changes (Levenshtein distance of less than 50%)
  - 190 other changes, about equal amounts function words and content words



- Focus on translation errors
  - not: basic spell checking
  - not: basic grammar checking
- Do not try the impossible
  - semantic errors
  - errors in function words
- What is left?
  - added content (insertions)
  - non-translated content (deletions)
  - inconsistency in terminology

- Word alignment of human translation and source
- Detect unaligned words
  - insertion of content words:  
unaligned sequence of words in the draft translation
  - deletion of content words:  
unaligned sequence of words in the source sentence
  - inconsistent terminology:  
source word occurs multiple times, aligned to different word
- Only content words (minimum 4 characters)

# Evaluation on Field Trial Data

- Two evaluation metrics
  - strict: predicted word X deleted / inserted
  - generous: predicted any deletion / insertion

Edit type	Strict Scoring		Generous Scoring		
	Precision	Recall	Precision	Recall	Baseline Precision
Deletion	7%	27%	11%	48%	7%
Insertion	-	-	5%	35%	4%
Any edit	-	-	20%	60%	14%

- Good enough to be useful?

# Subjective Evaluation

- Evaluation on community translation platform data
- English–German
- Predict insertions and deletions
- Manually check if these are valid suggestions (i.e., precision only) by native German speaker

# Results

- 4 cases of detection of valid errors (3 deletions, 1 insertion)
- 31 false alarms

Count	Type
16 cases	unaligned verb
6 cases	one-to-many alignment
2 cases	non-literal
6 cases	misalignment, often due to unknown word
1 case	valid verb ellipsis, repeated in sub clause

- Good enough to be useful?

# interactive translation prediction

## Input Sentence

Er hat seit Monaten geplant, im Oktober einen Vortrag in Miami zu halten.

## Professional Translator

|

## Input Sentence

Er hat seit Monaten geplant, im Oktober einen Vortrag in Miami zu halten.

## Professional Translator

| He



## Input Sentence

Er hat seit Monaten geplant, im Oktober einen Vortrag in Miami zu halten.

## Professional Translator

He | has

## Input Sentence

Er hat seit Monaten geplant, im Oktober einen Vortrag in Miami zu halten.

## Professional Translator

He has | for months

## Input Sentence

Er hat seit Monaten geplant, im Oktober einen Vortrag in Miami zu halten.

## Professional Translator

He planned |

## Input Sentence

Er hat seit Monaten geplant, im Oktober einen Vortrag in Miami zu halten.

## Professional Translator

He planned | for months

- Show  $n$  next words

Olvidarlo. Es demasiado | **arriesgado.** Estoy haciendo

- Show rest of sentence

# Spence Green's Lilt System

- Show alternate translation predictions

**C** Les étudiants eux-mêmes n'ont pas les moyens de se rendre à des cours, nous essayons de les aider de cette manière.

The students themselves cannot be required to attend courses, we are trying to help themselves cannot

**D** Dans le cadre de l'Institut Jedlička, nous transférerons ce projet dans un nouveau bâtiment.

themselves could not  
themselves do not  
themselves cannot afford

**E**

- Show alternate translations predictions with probabilities

To equip students with training, we have reduced mobility and Institute Jedlička, we will transfer this project to a new building.

■ routinely  
■ steadily  
■ regular  
■ regularly

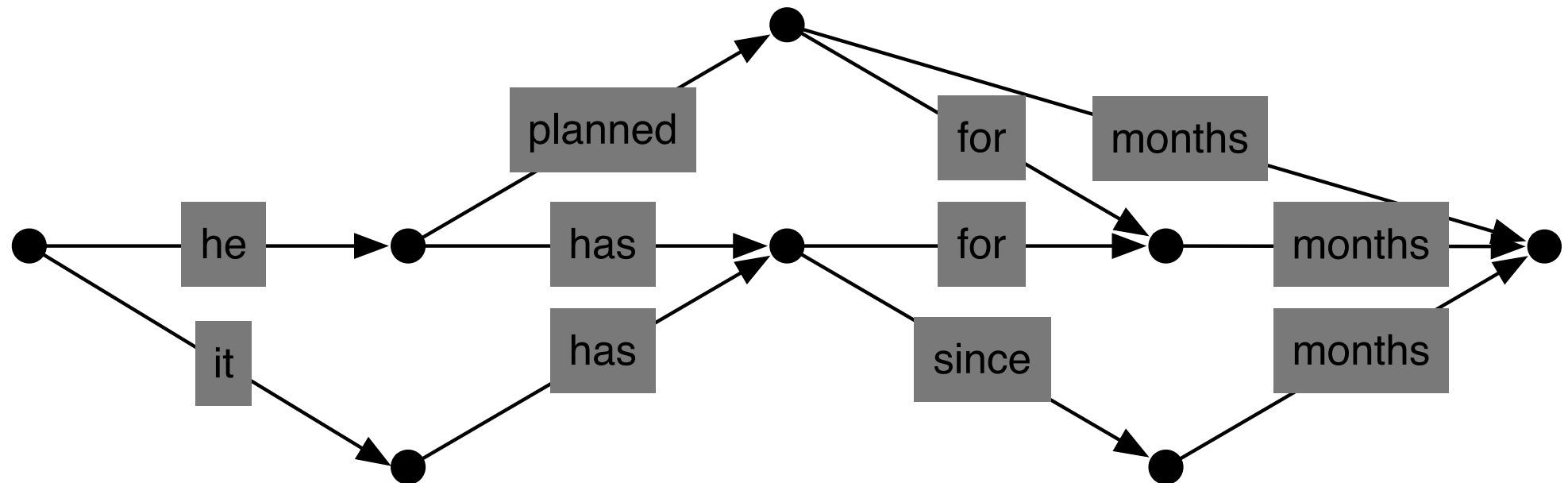
Des enseignants se rendent régulièrement auprès d'eux et proposent des activités qui les intéressent et les aident.

Teachers regularly visit Jedlička and propose activities that interest them and help them.

Les étudiants eux-mêmes n'ont pas les moyens de se rendre à des cours, nous essayons de les aider de cette manière.

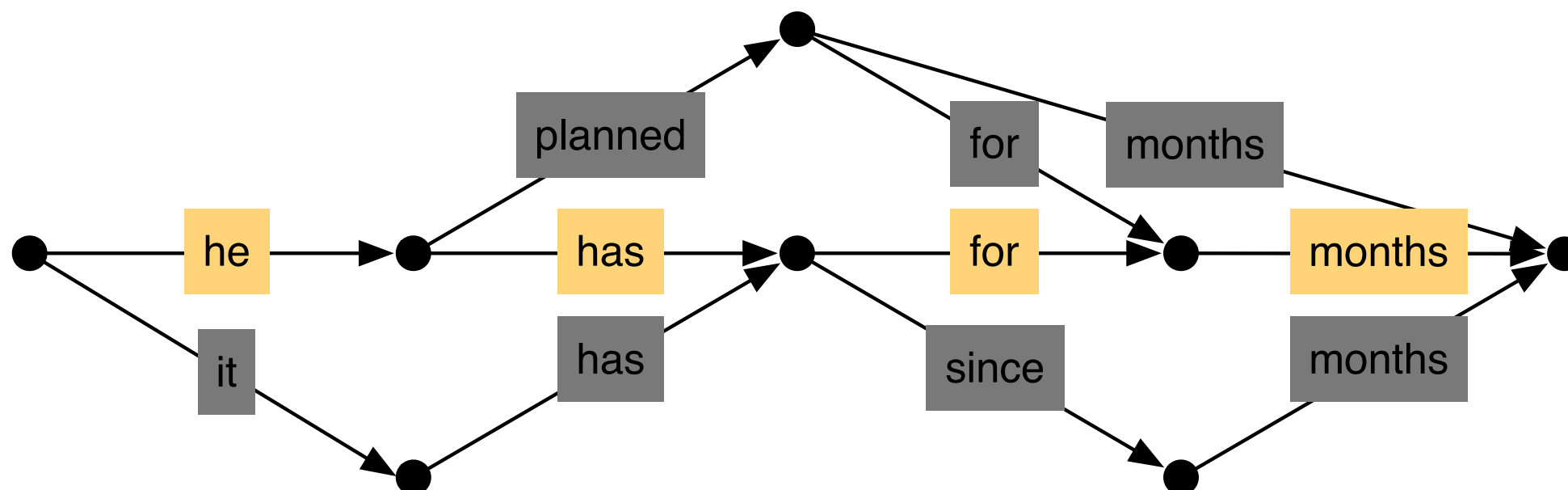
make regular  
are regularly

# Prediction from Search Graph



Search for best translation creates a graph of possible translations

# Prediction from Search Graph



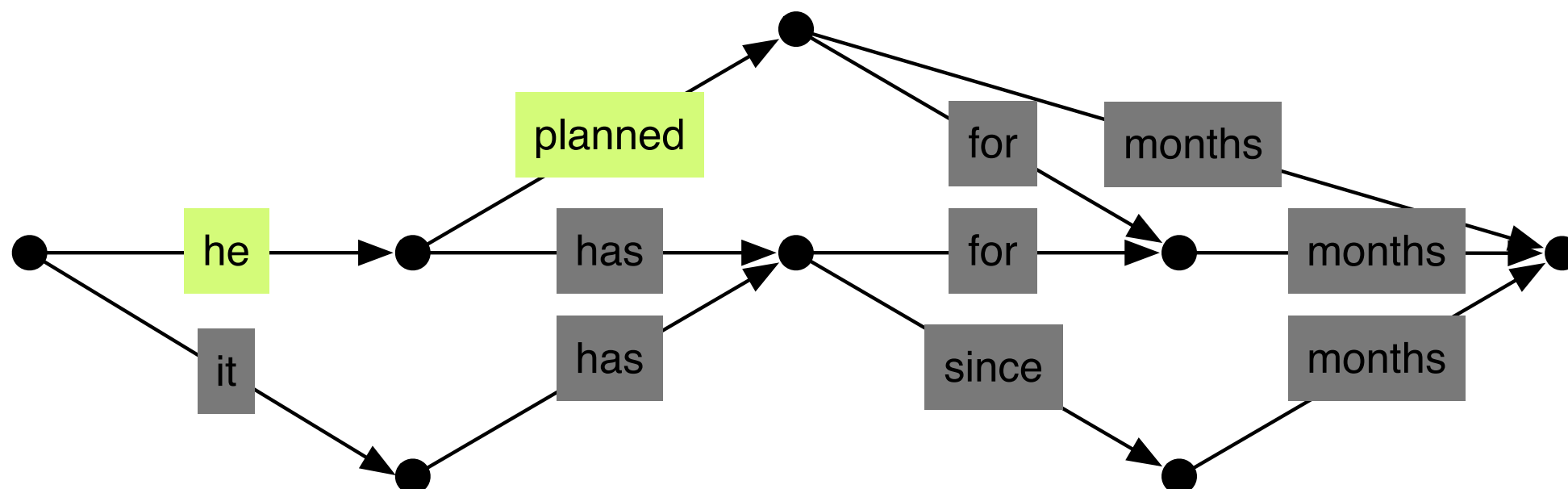
One path in the graph is the best (according to the model)

This path is suggested to the user



# Prediction from Search Graph

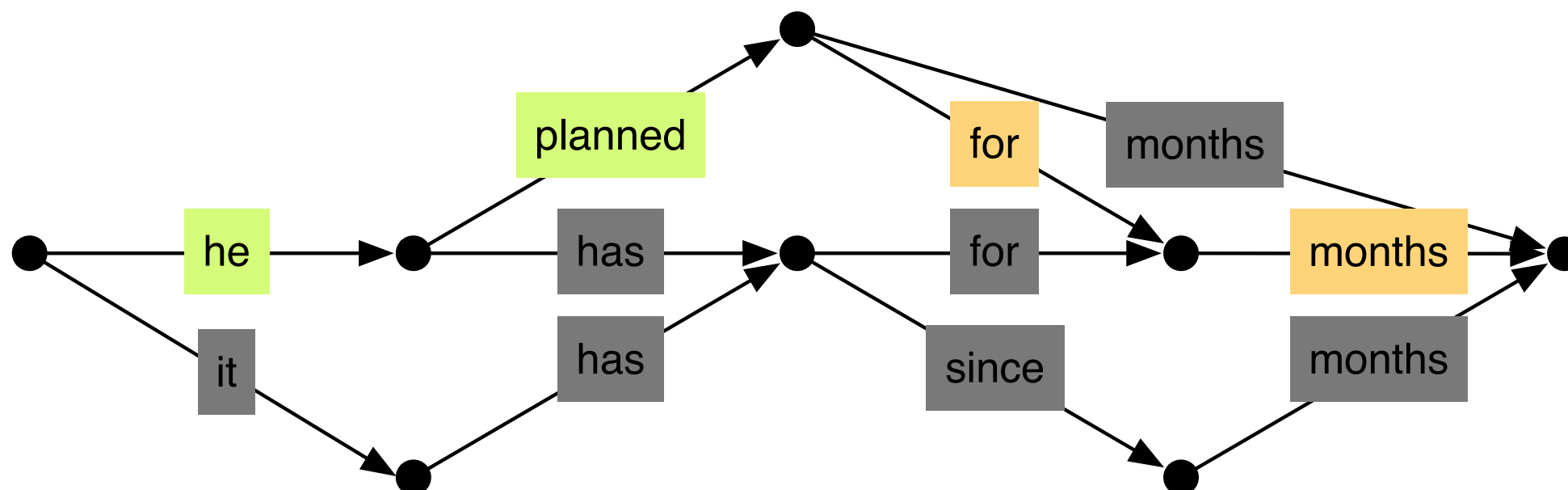
71



The user may enter a different translation for the first words

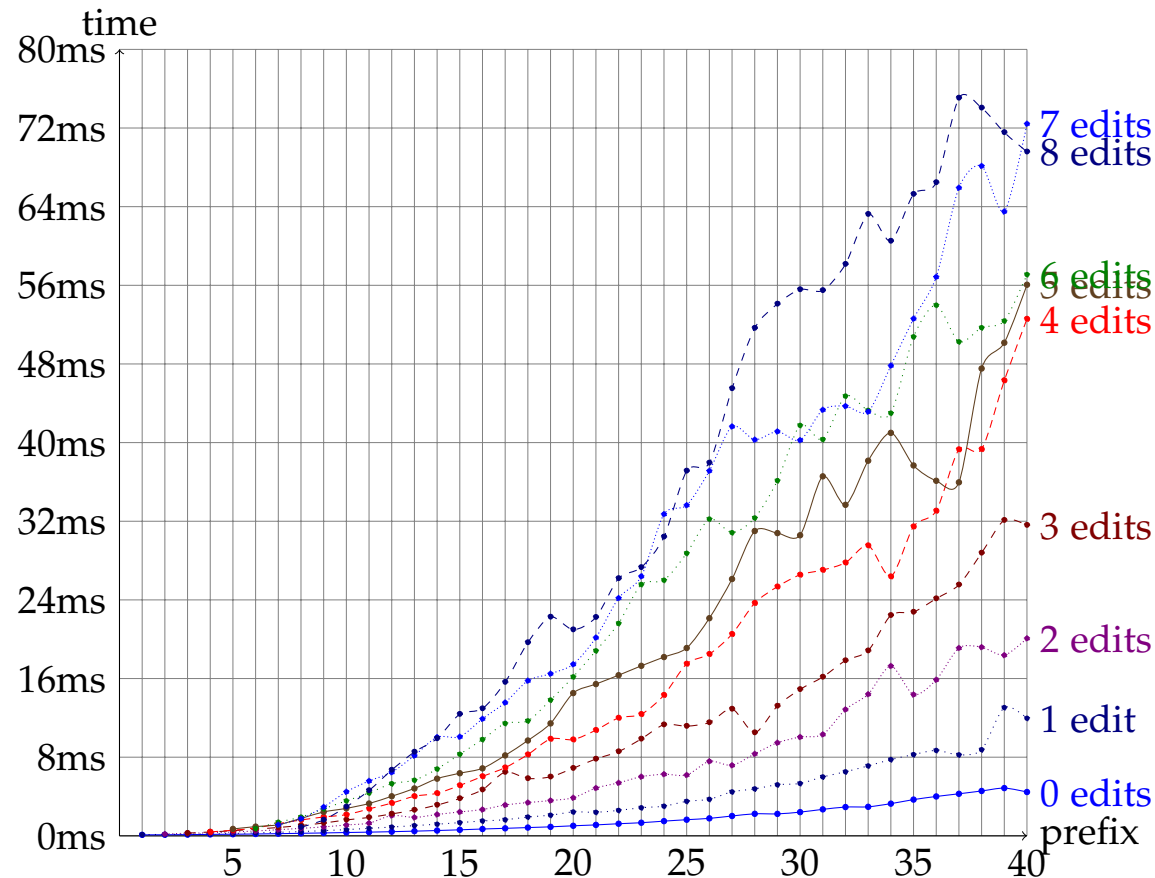
We have to find it in the graph

# Prediction from Search Graph



We can predict the optimal completion (according to the model)

# Speed of Algorithm



- Average response time based on length of the prefix and number of edits
- Main bottleneck is the string edit distance between prefix and path.

- Matching Last Word
  - more important to match last word in path
  - refinement of best path: search for last word
- Case-insensitive matching
- Approximate word matching
  - lower substitution cost for words that differ by a few letters
  - implemented at letter edit distance  $\leq 10\%$
- Stemmed matching
  - allow for difference in word endings (last 3 letters)
  - assumed to be morphological variation

# Word Completion

- Complete word once few letters are typed
- Example: predict *college* over *university*?
- User types the letter *u*  $\rightarrow$  change prediction
- "Desperate" word completion: find any word that matches

# Some Results

- News translation produced by post-editing MT output
- Same MT system used for simulated interactive translation prediction

#	Method	Word Acc.	Letter Acc.
1	Baseline	56.0%	75.2%
2	1 + Matching last word	59.0%	80.6%
3	2 + Case insensitive matching	58.7%	80.4%
4	2 + Approximate word matching	60.5%	80.6%
5	2 + Stemmed matching	59.4%	80.5%
6	4 + "Desperate" word completion	60.5%	84.5%

- Details see Koehn [ACL, 2014]

# Open Challenges

- Better metric than string edit distance to account for moves
- Retranslation or search graph matching?
- Interactive translation prediction for syntax-based models
  - syntax-based models work better for German, Chinese
  - search lattice → search forest
  - some preliminary work...
- Are neural machine translation models better at this?

⇒ Lots of interesting work in this area to be done



# bilingual concordancer



# Bilingual Concordancer

79



TRANSLATED

abandonner
×

abandon

ances des Etats-Unis à	<b>abandonner</b>	Musharraf -- et les co		merican reluctance to	<b>abandon</b>	Musharraf -- together
uridique, il a décidé d'	<b>abandonner</b>	la constitutionnalité, c		af has now decided to	<b>abandon</b>	constitutionality, remc
implement menacé d'	<b>abandonner</b>	ses accords commerci		simply threatened to	<b>abandon</b>	or never to conclude t

give up

erait donc contraint d'	<b>abandonner</b>	le droit de créer son p		e would be required to	<b>give up</b>	the right to develop it
n' était pas disposé à	<b>abandonner</b>	ses fonctions militaire		arraf was not ready to	<b>give up</b>	his military post, but a

to

t ne veulent donc pas	<b>abandonner</b>	leurs prérogatives dar		olicy and do not want	<b>to</b>	delegate this prerogat
-----------------------	-------------------	------------------------	--	-----------------------	-----------	------------------------

to abandon

es tout en refusant d'	<b>abandonner</b>	son arsenal nucléaire		drawal while refusing	<b>to abandon</b>	its nuclear weapons a
------------------------	-------------------	-----------------------	--	-----------------------	-------------------	-----------------------

# How does it Work?

- Have word-aligned parallel corpus
- Efficient data structure to quickly look up queried phrases (suffix arrays, we'll come back to them later)
- Translation spotting
  - look up queried phrase
  - use word alignment to identify target phrase
  - some edge cases (unaligned words at beginning/end)



# Linguee

English ↔ German ▼ ä ö ü ß

machine translation



Dictionary German-English

**machine translation** *noun*

**maschinelle Übersetzung** *f*

Maschinenübersetzung *f*

**translation machine** *noun*

**Übersetzungsmaschine** *f*

See also:

**machine** *n* — Maschine *f* · Gerät *nt* · Automat *m* · Anlage *f* · Apparat *m* · [...]

**machine** *v* — bearbeiten *vt* · maschinell herstellen *v* · spanen *v* · zerspanen *v* · maschinell bearbeiten *v* · [...]

**translation** *n* — Übersetzung *f* · Translation *f* · Übersetzen *nt* · Verschiebung *f* · Sprachübersetzung *f* · [...]

© Linguee Dictionary, 2015

Wikipedia

External sources (not reviewed)

The implementing provisions applicable to the machine translation system would have to be established by the Select Committee [...] [cep.eu](#)

Die Durchführungsbestimmungen für das System der maschinellen Übersetzung müssten vom engeren Ausschuss des EPO-Verwaltungsrats [...] [cep.eu](#)

By user licence agreements relating to the SYSTRAN machine translation software program concluded between the applicants' [...] [eur-lex.europa.eu](#)





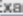



Durch Lizenzverträge über die Benutzung der Software für maschinelle Übersetzungen SYSTRAN zwischen den Rechtsvorgängern der Klägerinnen [...] [eur-lex.europa.eu](#)

[...] curriculum vitae, in forms suitable for multilingual machine translation, without restricting a user's option of adding other [...] [europarl.europa.eu](#)

[...] standardisierten Lebenslauf zu prüfen, die für eine automatische Übersetzung in mehrere Sprachen geeignet sind, wobei der Nutzer [...] [europarl.europa.eu](#)

# Verification of Terminology

- Translation of German *Windkraft*

<b>Examples</b>  	<b>Windkraft</b> (noun, feminine) (also: Windenergie)	 <b>wind power</b> (noun)	
	Zum Vergleich: <b>Windkraft</b> schafft fast sieben Mal mehr. ↳ German: <a href="http://www.goethe.de/wis/umw/thm/ntr/de92305.htm">www.goethe.de/wis/umw/thm/ntr/de92305.htm</a>	By way of comparison, <b>wind power</b> generates almost seven times as much. ↳ English: <a href="http://www.goethe.de/wis/umw/thm/ntr/en92305.htm">www.goethe.de/wis/umw/thm/ntr/en92305.htm</a>	
	Einführung von Windcube, einer neuen Generation von Wind Lidar für <b>Windkraft</b> . ↳ German: <a href="http://www.husumwindenergy.com/index.php?L...howUId]=1177">www.husumwindenergy.com/index.php?L...howUId]=1177</a>	Introducing Windcube, a new generation of <b>wind</b> Lidar for <b>wind power</b> . ↳ English: <a href="http://www.husumwindenergy.com/index.php?L...howUId]=1177">www.husumwindenergy.com/index.php?L...howUId]=1177</a>	
	<b>Windkraft</b> ist eine etablierte, wettbewerbsfähige Technologie mit hoher Zuverlässigkeit ↳ German: <a href="http://www.powergeneration.siemens.de/abo...ns-services/">www.powergeneration.siemens.de/abo...ns-services/</a>	<b>Wind power</b> is an established, competitive technology with high reliability ↳ English: <a href="http://www.powergeneration.siemens.com/abo...ns-services/">www.powergeneration.siemens.com/abo...ns-services/</a>	
<b>Examples</b>  	<b>Windkraft</b> (noun, feminine) (also: Windenergie)	 <b>wind energy</b> (noun)	
	Je mehr aber klimapolitische Sonntagsreden von der Politik auch in Taten umgesetzt werden, desto höher steigt dieser Preis und desto wettbewerbsfähiger werden saubere Energien wie die <b>Windkraft</b> . ↳ German: <a href="http://emagazine.credit-suisse.com/app/art...4382">emagazine.credit-suisse.com/app/art...4382</a> <=DE	But as the focus of the climate change issue shifts increasingly from policy to action, this price will increase and cleaner <b>energy</b> sources like <b>wind</b> will become more competitive. ↳ English: <a href="http://emagazine.credit-suisse.com/app/art...4382">emagazine.credit-suisse.com/app/art...4382</a> <=en	
	Nur wenige befürchten hingegen, dass dies auch bei erneuerbaren Energieträgern wie Biomasse oder <b>Windkraft</b> der Fall sein wird. ↳ German: <a href="http://www.eu2006.gv.at/de/News/Press_Rele...1proell.html">www.eu2006.gv.at/de/News/Press_Rele...1proell.html</a>	However, only a few fear that this will also be the case with renewable <b>energy</b> sources such as biomass or <b>wind energy</b> . ↳ English: <a href="http://www.eu2006.gv.at/en/News/Press_Rele...1proell.html">www.eu2006.gv.at/en/News/Press_Rele...1proell.html</a>	

- Context shows when each translation is used
- Indication of source supports trust in translations



UTILISATEUR : lapalme

REQUÊTES

MON COMPTE

PRÉFÉRENCES

AIDE

QUITTER

Signet / Favori personnalisé : TransSearch (qu'est-ce que c'est ?)

Requête bilingue

Collection de documents : Les Hansards canadiens

Expression : take+ .. ride

Chercher

92 traductions de **take+ .. ride** dans 106 occurrences

dindons de la farce	4
monté un bateau	3
faire avoir	3
se fasse rouler	2
fait berner	2
se fait jouer	2
moqués de	2
fait	2
les a	2
se sont fait avoir	2
le public pour attirer la	1
a fait une ballade	1
nous rouler dans ce projet nous tous	1
en train de monter un bateau à la population canadienne	1
tête des contribuables que se paie le	1
passer une petite vite	1
bourrer de l'autre côté de la chambre en	1
ont pris la voiture que pour faire une balade	1

## dindons de la farce

4

Emissions continue to rise and taxpayers are being **taken along for the ride**.

Les émissions continuent d'augmenter et c'est le contribuable qui est **le dindon de la farce**.

They are left with nothing. Now they are here illegally with no documentation. Canadians are being **taken for a ride**.

Ces personnes se trouvent ici illégalement, elles n'ont aucun document et nous, les Canadiens, sommes les **dindons de la farce**.

This would affect close to 400,000 Canadians, 80,000 of them Quebecers, who have been the ones **taken for a ride**.

Il s'agit d'une mesure qui toucherait près de 400 000 Canadiens, dont 80 000 Québécois, qui ont été les **dindons de la farce**.

I think that this is a prime example of a tainted system in which people who cannot afford to invest in sectors eligible for tax credits are urged to do so through all kinds of scams and end up being **taken for a ride**.

Je pense que c'est un exemple patent d'un système vicié, où des gens qui n'ont pas les moyens d'investir dans des domaines où on peut obtenir des crédits d'impôt se voient, par toutes sortes de subterfuges, invités à le faire et, en bout de ligne, ils se trouvent à être **les dindons de la farce**.

# TransSearch: Improved Transpotting

- Used to solve **difficult** translation problems
  - 7.2 million queries submitted to the system over a 6-year period
  - 87% contain at least two words
  - mainly search for idiomatic expressions such as *in keeping with*
- Improved translation spotting [Bourdaillet et al., MT Journal 2011]
- Filtering with classifier (45 features, trained on annotated data)
  - relative word count
  - word alignment scores
  - ratio of function words
- Merging of translations that only differ in function words, morphology
- Pseudo-relevance feedback



# translation options

# Translation Option Array

86



climbers are severely injured, and ten people are missing  
 after Mount Ontake (御嶽山, Ontake-san), a popular climbing  
 spot in central Japan, **erupted** for the first time in five years.

Kletterer sind schwer verletzt, und zehn Menschen werden  
 vermisst, nachdem Mount Ontake (御嶽山, Ontake-san), ein  
 beliebter Kletterplatz im zentralen Japan,

ausbruch, zum ersten

ITP    ≡    T→    DRAFT    **TRANSLATED**

Translation Options

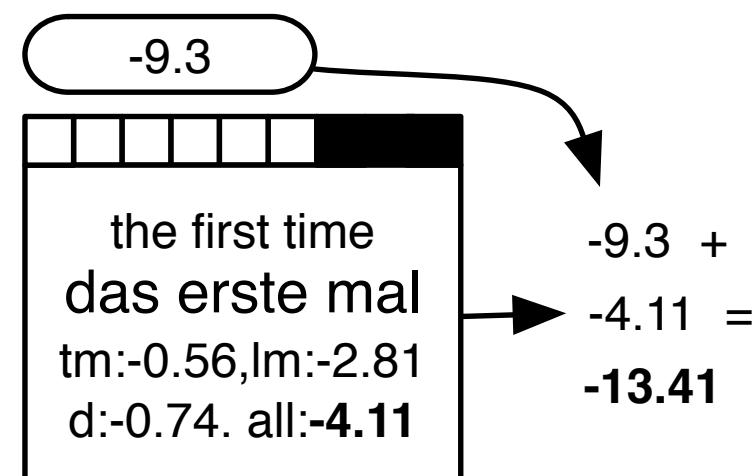
ke	-	san	)	,	a	popular	climbing	spot	in central	Japan	,	erupted	for the first time in five years	.
ke	-	san	)	,	ein	beliebtes	Klettern	vor Ort	in Mittel-	Japan,		ausbrach	zum ersten Mal in fünf Jahren	.
	und	San	)	,	ein	populär	Bergsteigen	vor	zentrale	Japan	,	ausbrach,	zum ersten Mal in	fünf Jahre.
	/		)	,	die	beliebt	Aufstieg	Fleck	zentralen	Japans,		platzte	zum ersten Mal	fünf Jahre
	der		)		eine	beliebte	abhalten,	ein, in	zentraler	Japan		Ausbruch		in fünf Jahren
	bis		)	,	in	populär	Erklimmen	Vor - Ort @-@	zentral	Japans	.	ausgebrochen	zum ersten Mal in der	von fünf Jahren.
	von		)	,	.	populär ist,	beim Besteigen	in	mittel-	in Japan	-	ausgebrochen ist	zum ersten Mal seit	fünf Jahren sind.

- Visual aid: non-intrusive provision of cues to the translator
- Trigger passive vocabulary



- Show up to 6 options per word or phrase
- Rank best option on top
- Use color highlighting to show likelihood  
(grey = less likely to be useful)
- Clickable: click on target phrase → added to edit area
- Automatic orientation
  - most relevant is next word to be translated
  - automatic centering on next word

- Basic idea: best options on top
- Problem: how to rank word translation vs. phrase translations?
- Method: utilize future cost estimates
- Translation score
  - sum of translation model costs
  - language model estimate
  - outside future cost estimate



# Improving Rankings

- Removal of duplicates and near duplicates

## bad

<b>erupted</b>
ausbrach
ausbrach,
platzte
Ausbruch
ausgebrochen
ausgebrochen ist

## good

<b>climbing</b>
Klettern
Bergsteigen
Aufstieg
abhalten,
Erklimmen
beim Besteigen

- Ranking by likelihood to be used in the translation  
→ can this be learned from user feedback?

# Enabling Monolingual Translators

- Monolingual translator
  - wants to understand a foreign document
  - has no knowledge of foreign language
  - uses a machine translation system■
- Questions
  - Is current MT output sufficient for understanding?
  - What else could be provided by a MT system?

- MT system output:

*The study also found that one of the genes **in the improvement in people with prostate cancer risk**, it also reduces the risk of suffering from diabetes.*

- What does this mean?■

- Monolingual translator:

*The research also found that one of the genes **increased people's risk of prostate cancer, but at the same time lowered people's risk of diabetes.**■*

- Document context helps

# Example: Arabic

92



وكان	مجلس	النواب	الاميركي	اعتمد	الخميس	قانونا	يطالب	يسحب	القوات	المقاتلة	الاميركية	من	العراق	في	موعد	اقصاه	الاول	من	نيسان	@/@@	ابريل
the	the us house of representatives	adopted	thursday	legally	calls for the withdrawal of	combat troops	us	iraq	in	no later than	the first	from	april								
the us house of representatives	the	thursday ,	law		the fighting forces	the us	from iraq		the latest	the first of	april										
the us house	adopted the	thu	the legally		fighting forces	us	from iraq in			i	april										
it was	us house of representatives	was adopted	thursday , the	the law	demands withdrawal of troops	fighter	the us		no later than	first	on april										
he was	the us house	adopted by	thursday 's	a law	calls for withdrawal of	combat forces	of		in the	not later than	first of										
he	us house	adopted by the	on thursday	a legally	calls for the withdrawal	forces	the fighter														
earlier ,		us	adopted a	on thursday ,	by law	demands the withdrawal of	troops														
was			, was adopted	thursday the	legally ,	demands withdrawal of															
it was the			adopted ,	thu ,	the legal	calls for withdrawal															
earlier , the			adopted , the	thursday , a	legally @-@	demands the withdrawal															
2008 ,	متحديا	مرة	جديدة	الرئيس	جورج	بوش	الذي	يعارض	اي	تحديد	موعد										
2008 ,	defying	once	new	president george w. bush	which opposes the	no date has been set for the															
the 2008	defiant	once again		president george bush	who opposes	no date has been set for															
2008	challenging	again	the new		, which opposes	no date has been set															
	a defiant	the first			, who opposes the	a date .															
	in defiance of	once again ,			, who opposes	date .															
	, challenging	once again the		president george bush , who	opposed to setting any	the date of the															
	, in defiance	for the first time	a new	president george w. bush 's	which opposes	no date															
in 2008 ,	defying the	again		us president george w. bush	opposed to	any	the date of														
	challenging the	time			who opposes the	date of															
	, defying	once again , the			opposes	date															

up to 10 translations for each word / phrase

# Example: Arabic

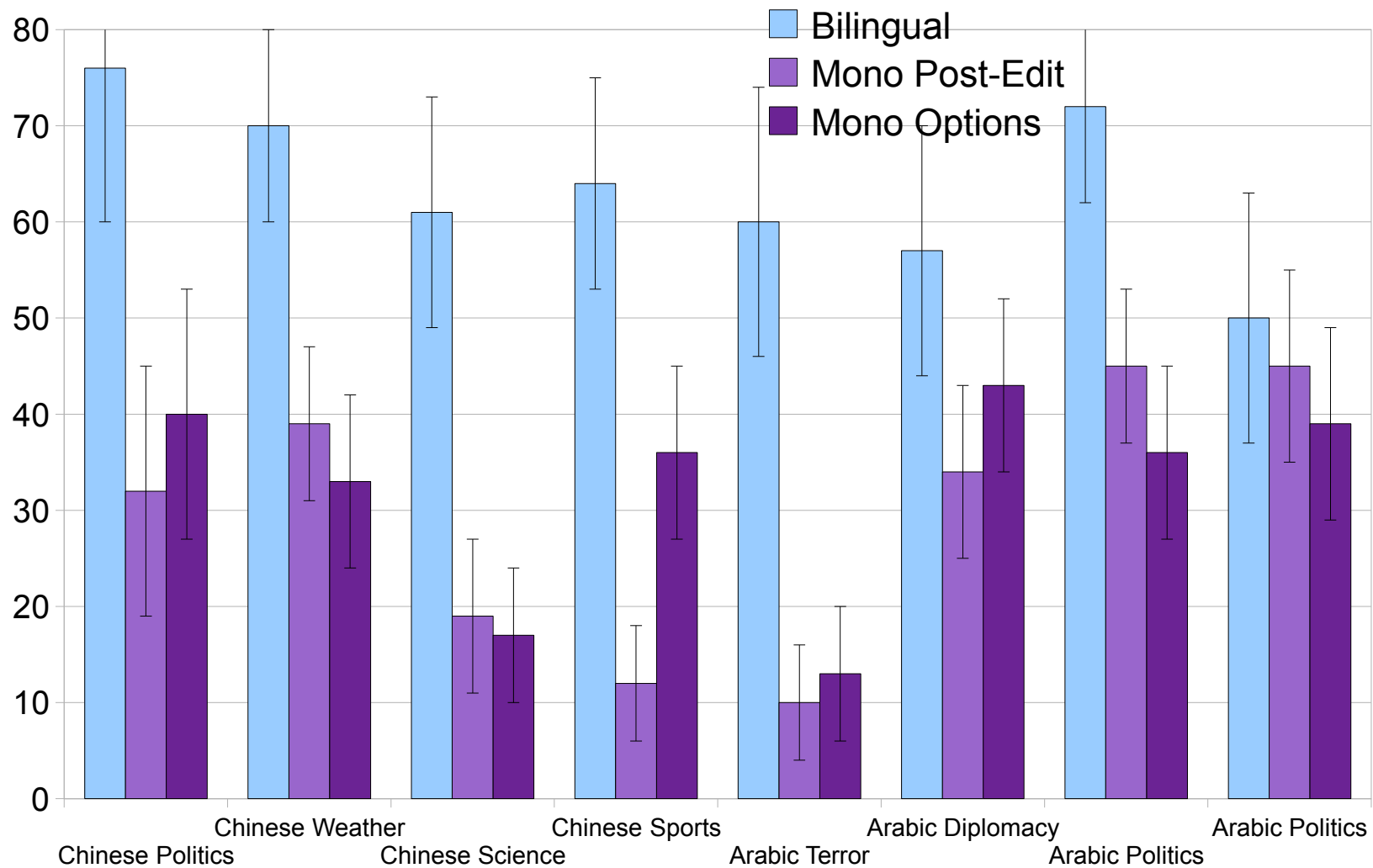
93



العراق	من	الاميركية	المقاتلة	القوات	يسحب
iraq	us	the us	fighting forces	combat troops	withdrawal of
from iraq	us	the us	fighting forces	combat troops	the fighting forces
from irac	us	the us	fighter	combat forces	withdrawal of troops
i	of	the us	the fighter	troops	withdrawal of
	from	the us			the withdrawal
iraq	of the	the us			the withdrawal of
ir	from iraq	the us			the withdrawal
	the american	the us			the withdrawal

# Monolingual Translation with Options

94



No big difference — once significantly better



# Monolingual Translation Triage

- Study on Russian–English (Schwartz, 2014)
- Allow monolingual translators to assess their translation
  - confident → accept the translation
  - verify → proofread by bilingual
  - partially unsure → part of translation handled by bilingual
  - completely unsure → handled by bilingual
- Monolingual translator highly effective in triage



- Main findings
  - monolingual translators may be as good as bilinguals■
  - widely different performance by translator / story■
  - named entity translation critically important■
- Various human factors important
  - domain knowledge■
  - language skills■
  - effort



# paraphrasing

## Input Sentence

Er hat seit Monaten geplant, im Oktober einen Vortrag in Miami zu halten.

## Professional Translator

He planned for months to give a lecture in Miami in October.

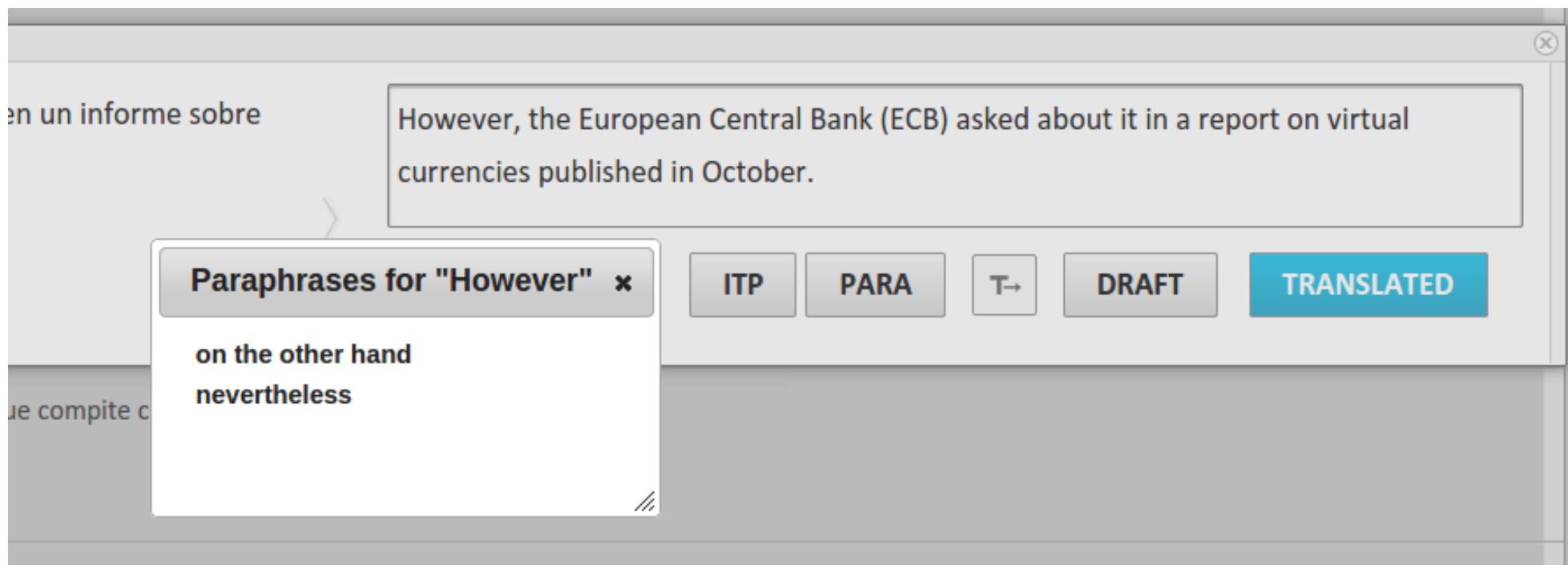
give a presentation

present his work

give a speech

speak

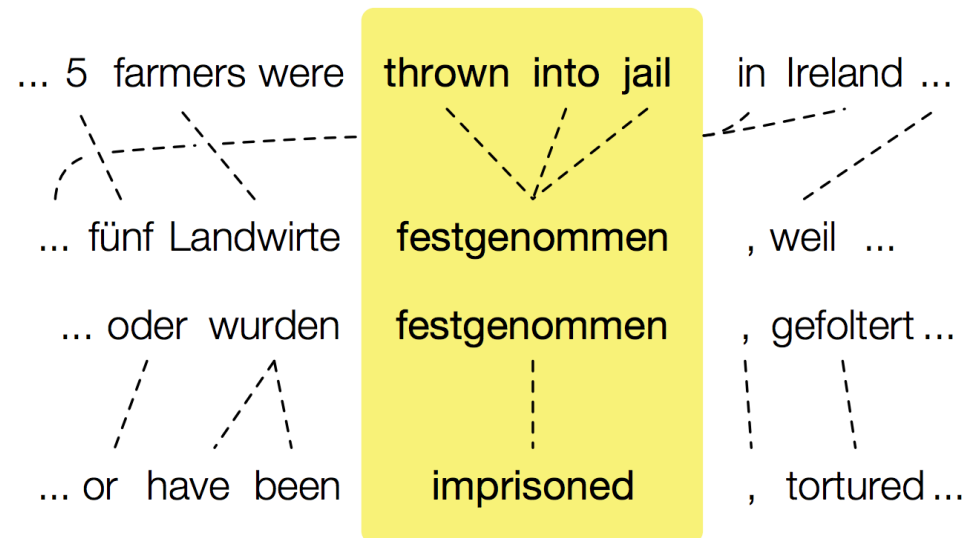
User requests alternative translations for parts of sentence.



- User marks part of translation
- Clicks on paraphrasing button
- Alternative translations appear

- Somewhat popular research area
- Popular method: extract from parallel data

- goal: find paraphrases for phrase  $e$
  - look up likely translations  $f_1, f_2, \dots$  for  $e$
  - for each  $f_i$ , look up likely translations  $e'_{i1}, e'_{i2}, \dots$
- ⇒ these are the paraphrases



- Refinement: collect over several foreign languages, intersect
- Paraphrase database for several languages:  
<http://paraphrase.org/>



- Our problem: paraphrasing in context
    - driven by source
    - considers sentence context
    - ranking and diversity important
    - real time performance
  - Approach
    - target span is mapped to source span
    - search graph is consulted for alternative translations for source span
    - additional translations generated by combining translation options
- ⇒ initial list of translations
- various components to distill  $n$ -best paraphrases



- Filtering: remove some translations
  - with extraneous punctuation
  - too similar to others
  - additional function words
- Scoring: score translations
  - translation model scores
  - language model score in context
  - compare alternate translations against best path
- Sorting: rank list
  - cluster translations by similarity
  - picks best translation from each cluster





- Motivation
  - alternative translations should fix translation errors
  - create bad translations by back-translation
- Process
  - Train machine translation system for both directions
  - Translate test set  $\text{target} \rightarrow \text{source} \rightarrow \text{target}^*$
  - Spot differences between  $\text{target}$  and  $\text{target}^*$
  - Use span in  $\text{target}^*$  as “marked by user”, span in  $\text{target}$  as correct

# Example



- Translate

*Unlike in Canada , **the American states**  
are responsible for the organisation of federal  
elections.*

- Into

**В отличие от Канады, американские штаты  
ответственны за организацию федеральных  
выборов в соединенных штатах .**

- Back into English

*Unlike in Canada , **US states**  
are responsible for the organization of  
federal elections.*

- Web based interactive evaluation tool
- Same setup as automatic evaluation
  - shows target span
  - 5 selectable paraphrases
  - user accepts one → correct
- Four users (U1–U4)
- Number of instances where one translation is correct

Method	U1	U2	U3	U4	average score
1	8	6	9	6	6/50
7	15	17	12	10	13/50
10	24	20	26	29	26/50

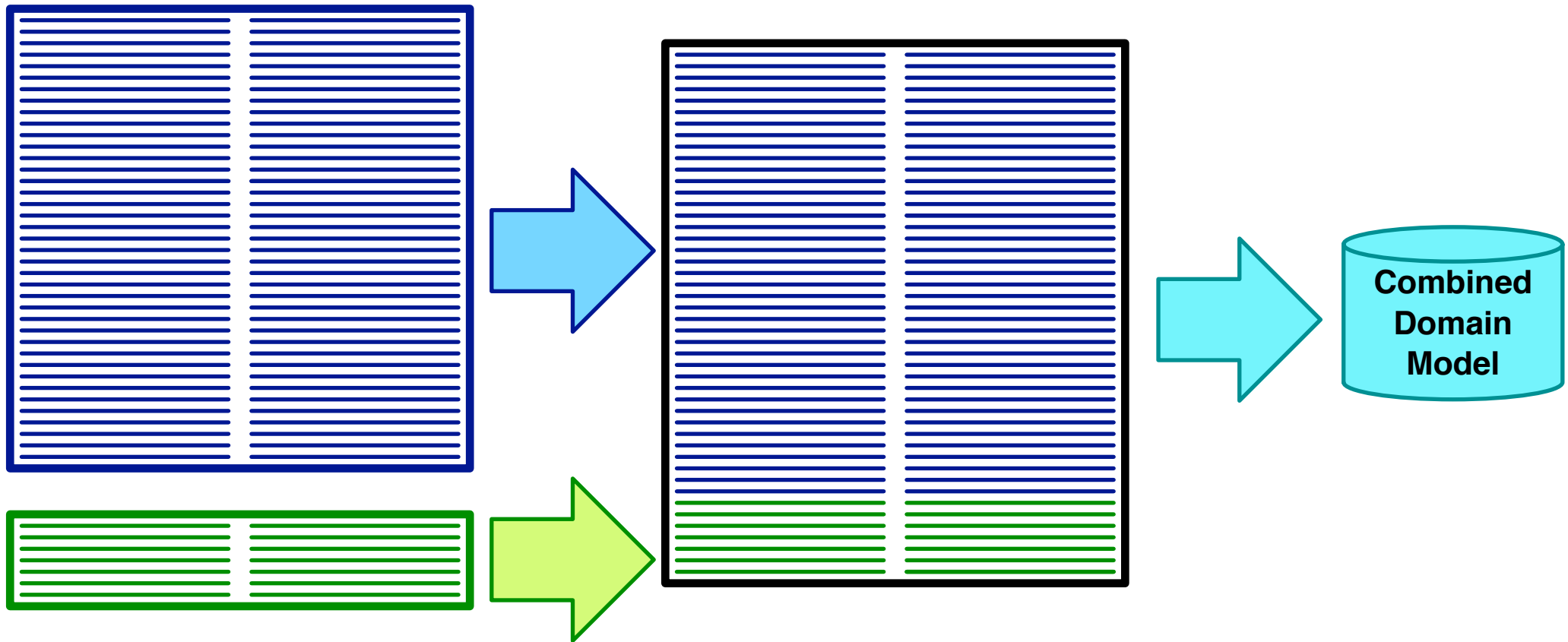


# adaptation

- Machine translation works best if optimized for domain
- Typically, large amounts of out-of-domain data available
  - European Parliament, United Nations
  - unspecified data crawled from the web
- Little in-domain data (maybe 1% of total)
  - information technology data
  - more specific: IBM's user manuals
  - even more specific: IBM's user manual for same product line from last year
  - and even more specific: sentence pairs from current project
- Various domain adaptation techniques researched and used

# Combining Data

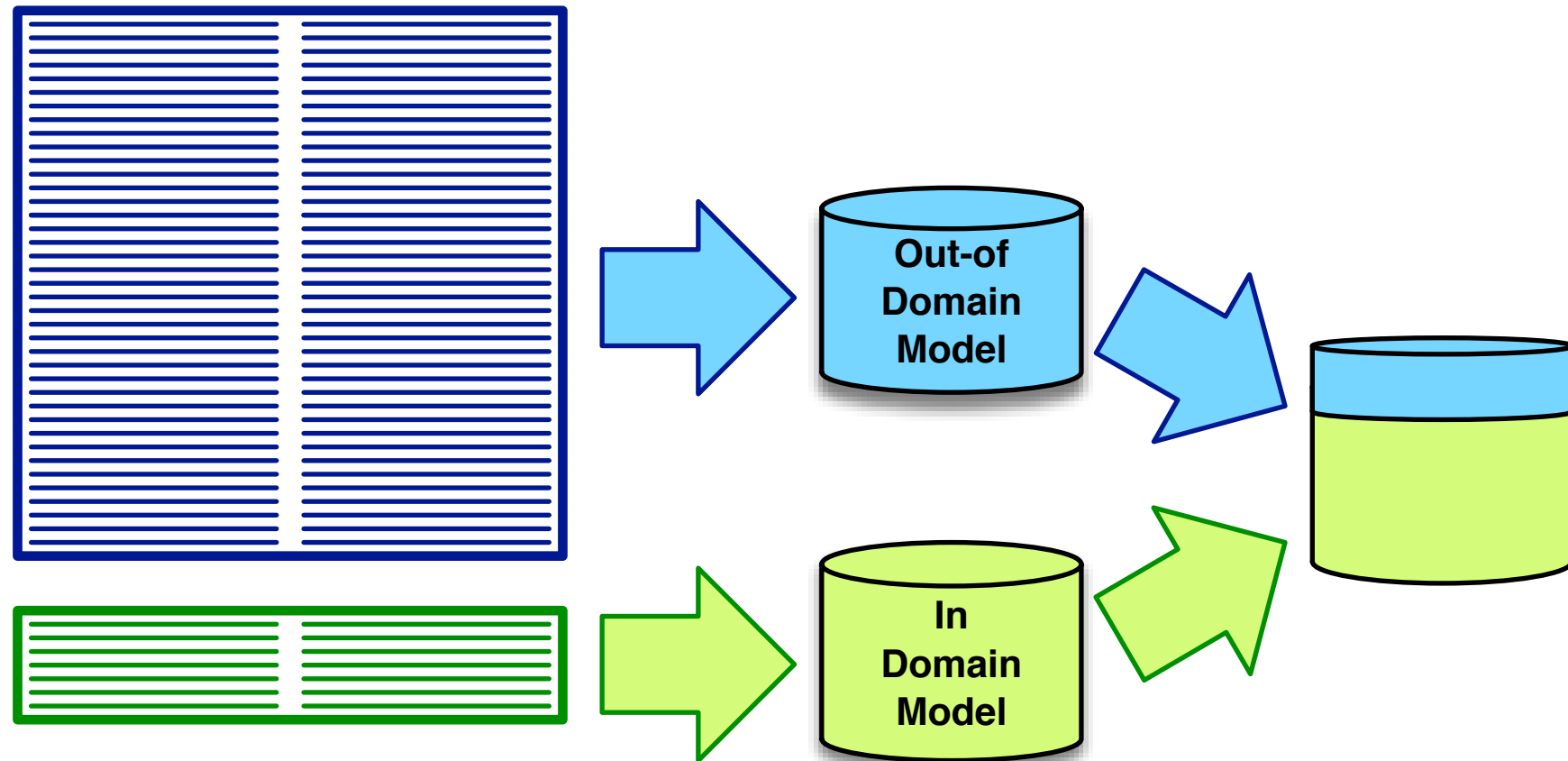
108



- Too biased towards out of domain data
- May flag translation options with indicator feature functions

# Interpolate Models

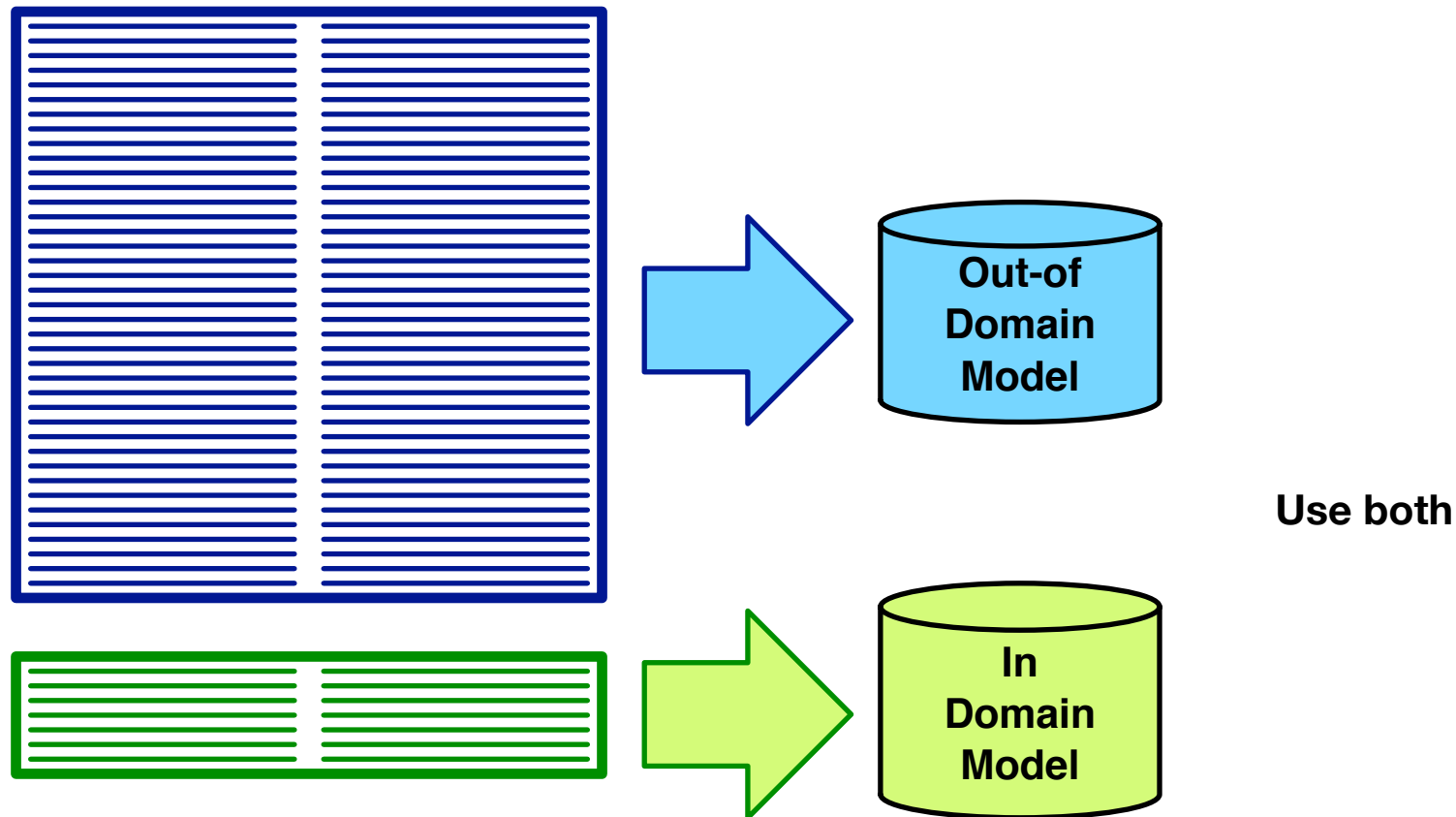
109



- $p_c(e|f) = \lambda_{\text{in}}p_{\text{in}}(e|f) + \lambda_{\text{out}}p_{\text{out}}(e|f)$
- Quite successful for language modelling

# Multiple Models

110

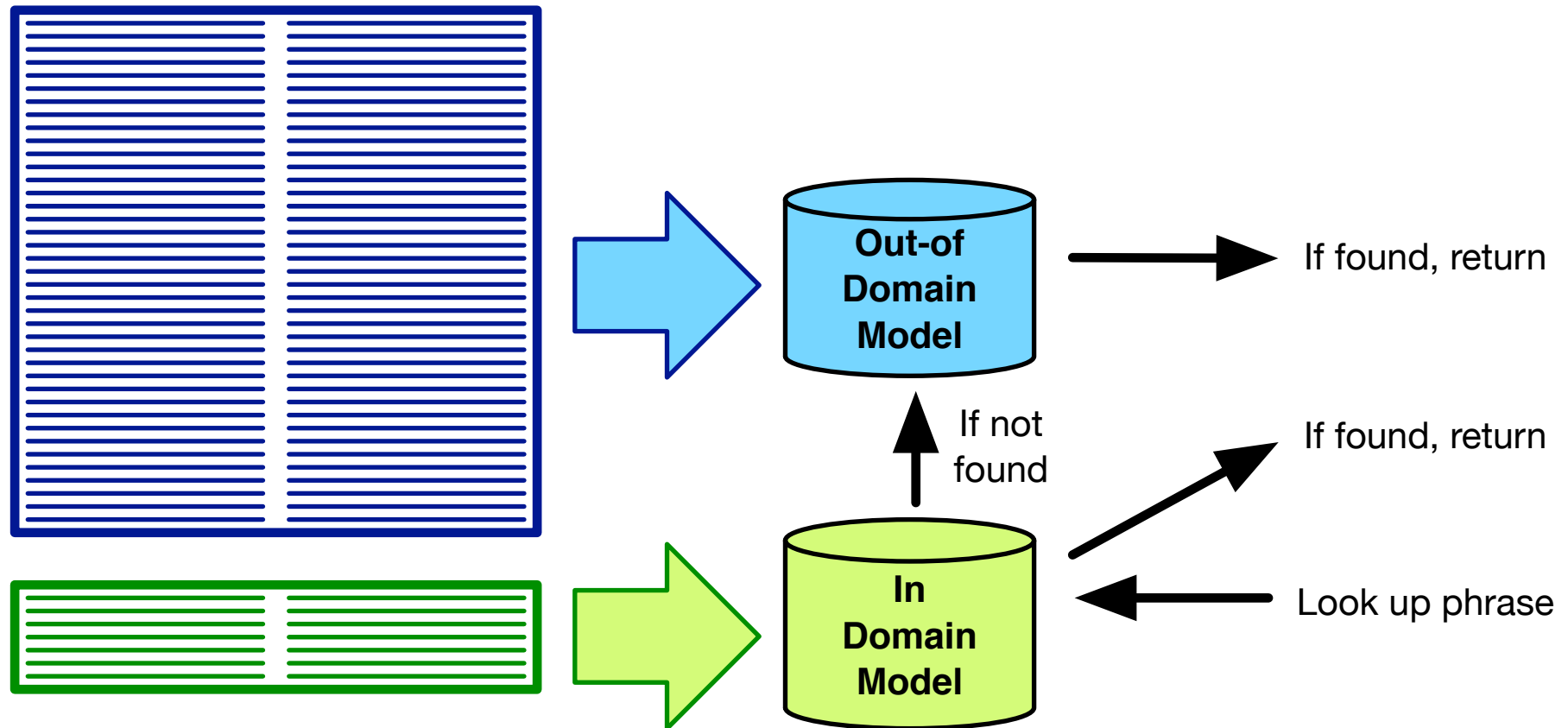


- Multiple models  $\rightarrow$  multiple feature functions



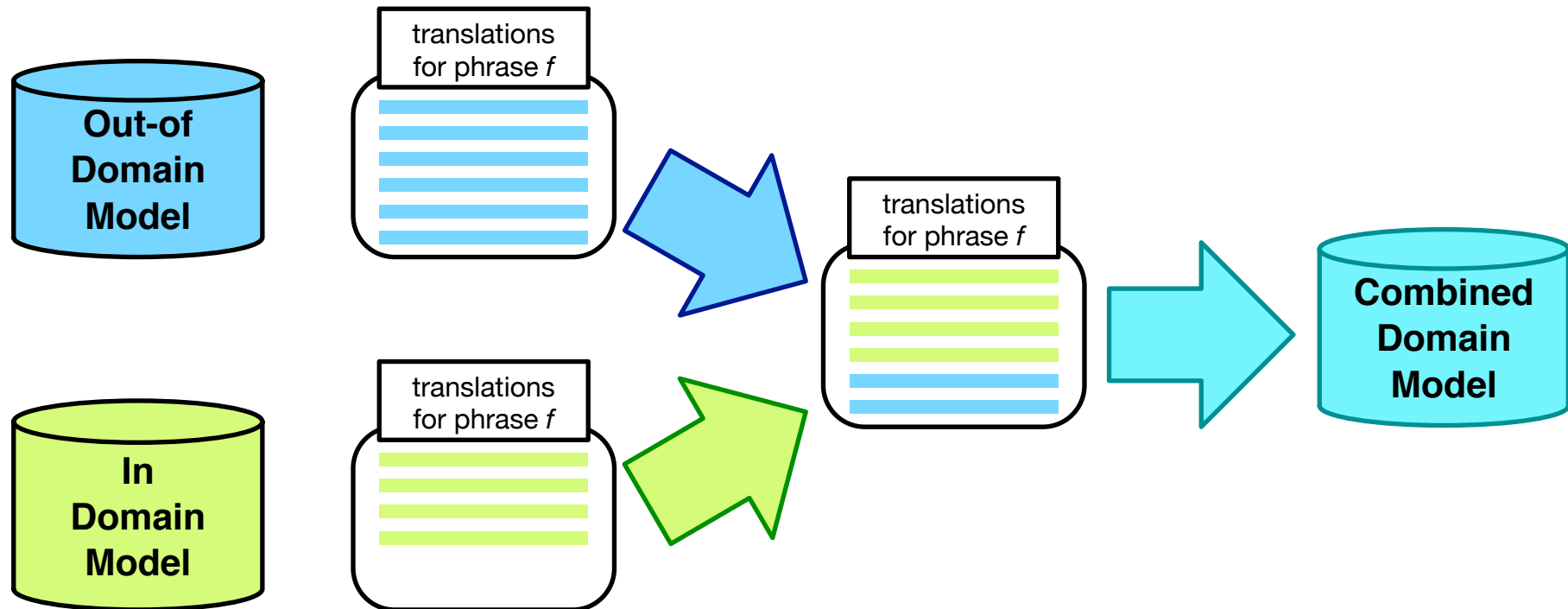
# Backoff

111



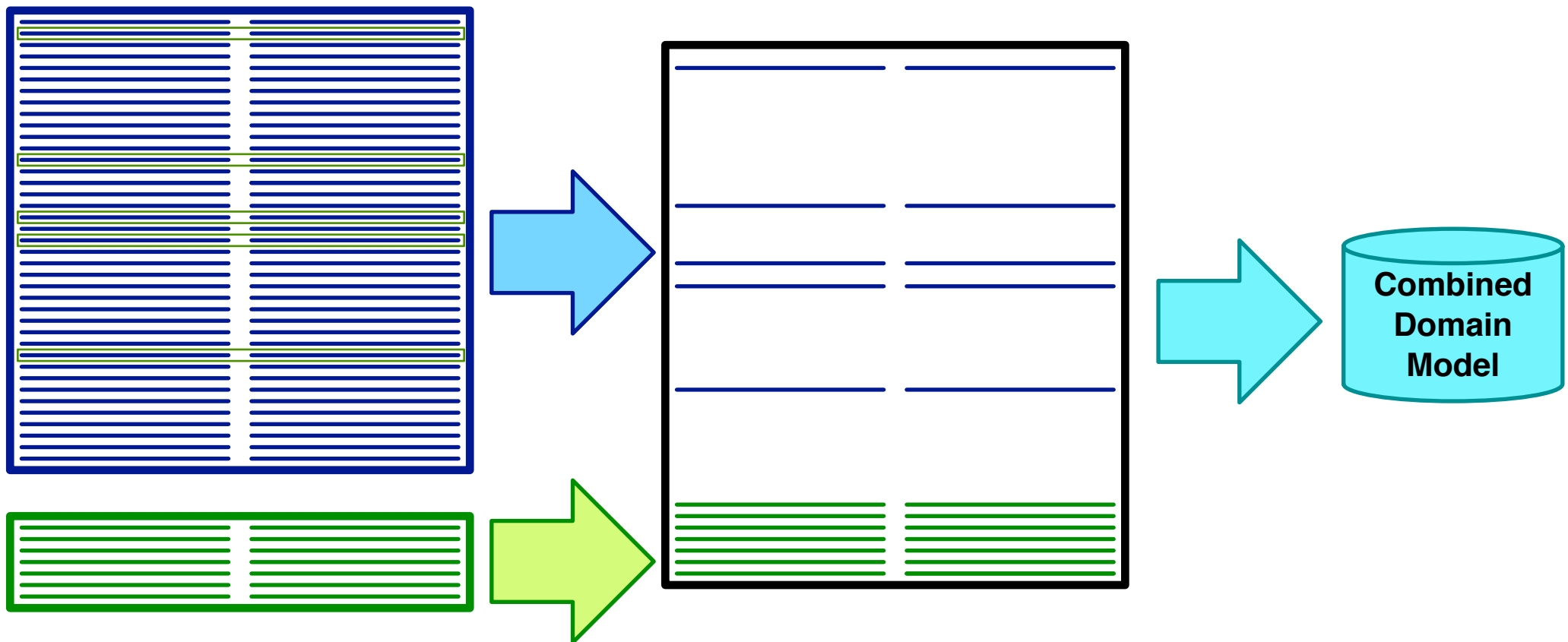
# Fill-Up

112



- Use translation options from in-domain table
- Fill up with additional options from out-of-domain table

# Sentence Selection

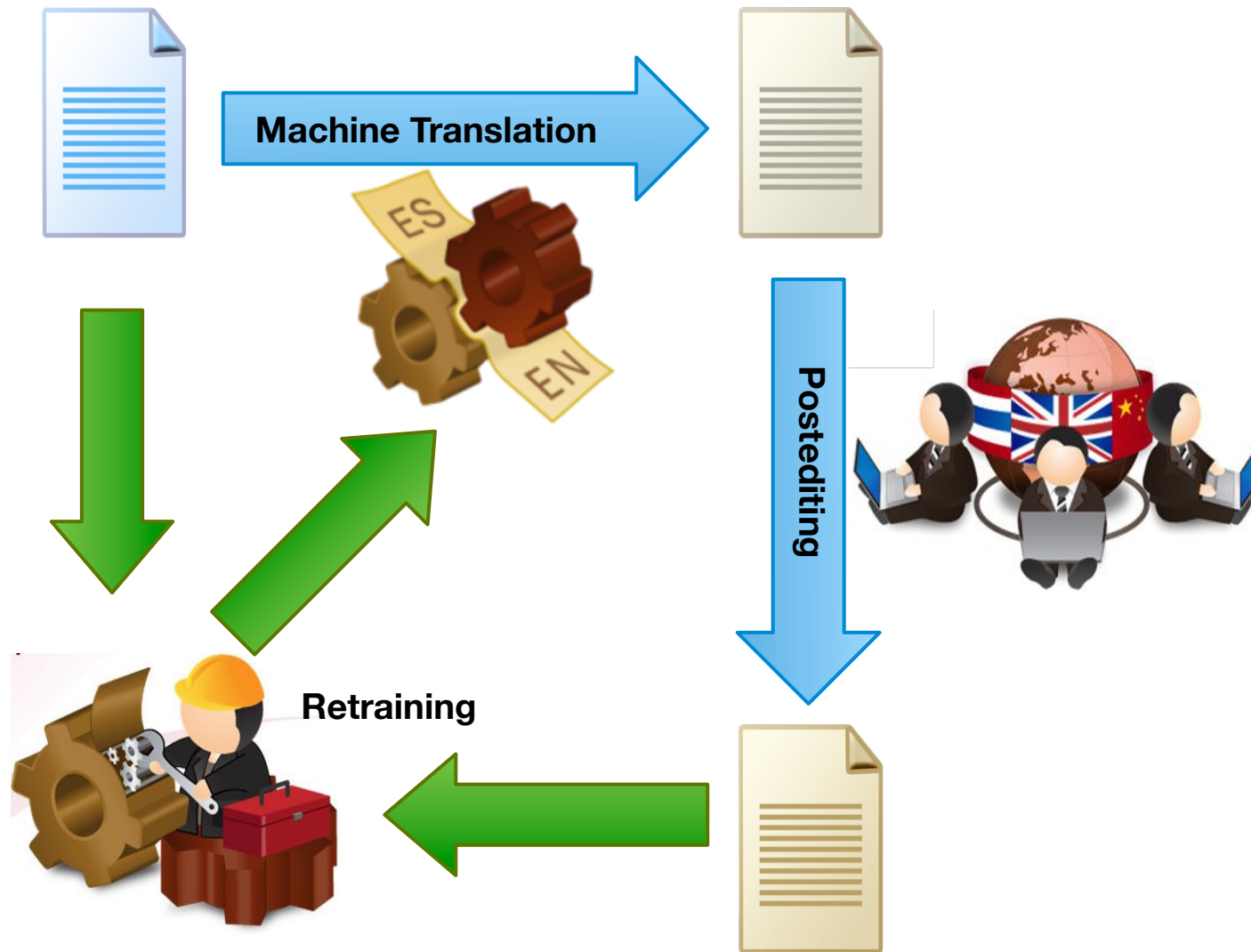


- Select out-of-domain sentence pairs that are similar to in-domain data
- Score similarity with language model, other means

- Method developed by the Matecat project
- Update model during translation project
- After each day
  - collected translated sentences
  - add to model
  - optimize
- Main benefit after the first day

# Instant Adaptation

115



# Adaptable Translation Model

116



- Store in memory
  - parallel corpus
  - word alignment
- Adding new sentence pair
  - word alignment of sentence pair
  - add sentence pair
  - update index (suffix array)
- Retrieve phrase translations on demand

- Needed: word alignment method that scores a sentence pairs
- Online EM algorithm
  - keep sufficient statistics of corpus in memory
  - run EM iteration on single sentence pair
  - update statistics
  - return word alignment
- For efficiency reason, a static model may be sufficient
- Implementations in both mGIZA and fast-align

# Suffixes

118



- 1 government of the people , by the people , for the people
- 2 of the people , by the people , for the people
- 3 the people , by the people , for the people
- 4 people , by the people , for the people
- 5 , by the people , for the people
- 6 by the people , for the people
- 7 the people , for the people
- 8 people , for the people
- 9 , for the people
- 10 for the people
- 11 the people
- 12 people



# Sorted Suffixes

119



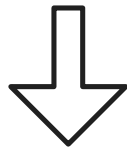
5 , by the people , for the people  
9 , for the people  
6 by the people , for the people  
10 for the people  
1 government of the people , by the people , for the people  
2 of the people , by the people , for the people  
12 people  
4 people , by the people , for the people  
8 people , for the people  
11 the people  
3 the people , by the people , for the people  
7 the people , for the people

# Suffix Array

120



5	, by the people , for the people
9	, for the people
6	by the people , for the people
10	for the people
1	government of the people , by the people , for the people
2	of the people , by the people , for the people
12	people
4	people , by the people , for the people
8	people , for the people
11	the people
3	the people , by the people , for the people
7	the people , for the people



suffix array: sorted index of corpus positions

# Querying the Suffix Array

121



5	, by the people , for the people
9	, for the people
6	by the people , for the people
10	for the people
1	government of the people , by the people , for the people
2	of the people , by the people , for the people
12	people
4	people , by the people , for the people
8	people , for the people
11	the people
3	the people , by the people , for the people
7	the people , for the people

Query: people

# Querying the Suffix Array

122



5	, by the people , for the people
9	, for the people
6	by the people , for the people
10	for the people
1	government of the people , by the people , for the people
2	of the people , by the people , for the people
12	people
4	people , by the people , for the people
8	people , for the people
11	the people
3	the people , by the people , for the people
7	the people , for the people

Query: **people**

Binary search: start in the middle

# Querying the Suffix Array

123



5	, by the people , for the people
9	, for the people
6	by the people , for the people
10	for the people
1	government of the people , by the people , for the people
2	→ of the people , by the people , for the people
12	people
4	people , by the people , for the people
8	people , for the people
11	the people
3	the people , by the people , for the people
7	the people , for the people

Query: **people**

Binary search: discard upper half

# Querying the Suffix Array

124



5	, by the people , for the people
9	, for the people
6	by the people , for the people
10	for the people
1	government of the people , by the people , for the people
2	of the people , by the people , for the people
12	people
4	people , by the people , for the people
8	people , for the people
11	the people
3	the people , by the people , for the people
7	the people , for the people

Query: **people**

Binary search: middle of remaining space

# Querying the Suffix Array

125



5	, by the people , for the people
9	, for the people
6	by the people , for the people
10	for the people
1	government of the people , by the people , for the people
2	of the people , by the people , for the people
12	people
4	people , by the people , for the people
8	people , for the people
11	the people
3	the people , by the people , for the people
7	the people , for the people

Query: people

Binary search: match

# Querying the Suffix Array

126



5		, by the people , for the people
9		, for the people
6		by the people , for the people
10		for the people
1		government of the people , by the people , for the people
2	→	of the people , by the people , for the people
12	↙	people
4	↘	people , by the people , for the people
8	↘	people , for the people
11	↙	the people
3	↘	the people , by the people , for the people
7		the people , for the people

Query: **people**

Finding matching range with additional binary searches for start and end



- Cache-based models
- Language model
  - give bonus to n-grams in previous user translation
- Translation model
  - give bonus to translation options in previous user translation
- Decaying score for bonus (less recent, less relevant)



# integration of translation memories



- **Translation Memory (TM)**

- translators store past translation in database
- when translating new text, consult database for similar segments
- fuzzy match score defines similarity

widely used by translation agencies■

- **Statistical Machine Translation (SMT)**

- collect large quantities of translated text
- extract automatically probabilistic translation rules
- when translating new text, find most probable translation given rules

wide use of free web-based services  
not yet used by many translation agencies

# TM

# vs.

# SMT

130



used by  
human translator

used by  
target language information seeker

restricted domain  
(e.g. product manual)

open domain translation  
(e.g. news)

very repetitive content

huge diversity (esp. web)

corpus size:  
1 million words

corpus size:  
100-1000 million words

commercial developers  
(e.g., SDL Trados)

academic/commercial research  
(e.g., Google)



- Input

The second paragraph of Article 21 is deleted .

- Fuzzy match in translation memory

The second paragraph of Article 5 is deleted .

⇒ **Part of the translation from TM fuzzy match**

**Part of the translation with SMT**

The second paragraph of Article 21 is deleted .

# Example

132



- Input sentence:

The second paragraph of Article 21 is deleted .

# Example

133



- Input sentence:

The second paragraph of Article 21 is deleted .

- Fuzzy match in translation memory:

The second paragraph of Article 5 is deleted .

=

À l' article 5 , le texte du deuxième alinéa est supprimé .

# Example

134



- Input sentence:

The second paragraph of Article 21 is deleted .

- Fuzzy match in translation memory:

The second paragraph of Article 5 is deleted .

=

À l' article 5 , le texte du deuxième alinéa est supprimé .

- Detect mismatch (string edit distance)



# Example

135



- Input sentence:

The second paragraph of Article 21 is deleted .

- Fuzzy match in translation memory:

The second paragraph of Article 5 is deleted .

=

À l' article 5 , le texte du deuxième alinéa est supprimé .

- Detect mismatch (string edit distance)
- Align mismatch (using word alignment from GIZA++)

# Example



- Input sentence:

The second paragraph of Article 21 is deleted .

- Fuzzy match in translation memory:

The second paragraph of Article 5 is deleted .

=

À l' article 5 , le texte du deuxième alinéa est supprimé .

Output word(s) taken from the target TM

# Example

137



- Input sentence:

The second paragraph of Article 21 is deleted .

- Fuzzy match in translation memory:

The second paragraph of Article 5 is deleted .

=

À l' article 5 , le texte du deuxième alinéa est supprimé .

Output word(s) taken from the target TM

Input word(s) that still need to be translated by SMT

# Example

138



- Input sentence:

The second paragraph of Article 21 is deleted .

- Fuzzy match in translation memory:

The second paragraph of Article 5 is deleted .

=

À l' article 5 , le texte du deuxième alinéa est supprimé .

- XML frame (input to Moses)

<xml translation=" À l' article " /> 21

<xml translation=" , le texte du deuxième alinéa est supprimé . " />

# Example

139



- Input sentence:

The second paragraph of Article 21 is deleted .

- Fuzzy match in translation memory:

The second paragraph of Article 5 is deleted .

=

À l' article 5 , le texte du deuxième alinéa est supprimé .

- More compact formalism for the purposes of this presentation:

< À l' article > 21 < , le texte du deuxième alinéa est supprimé . >



- XML frames

<À l' article> 21 <, le texte du deuxième alinéa est supprimé .>

for input

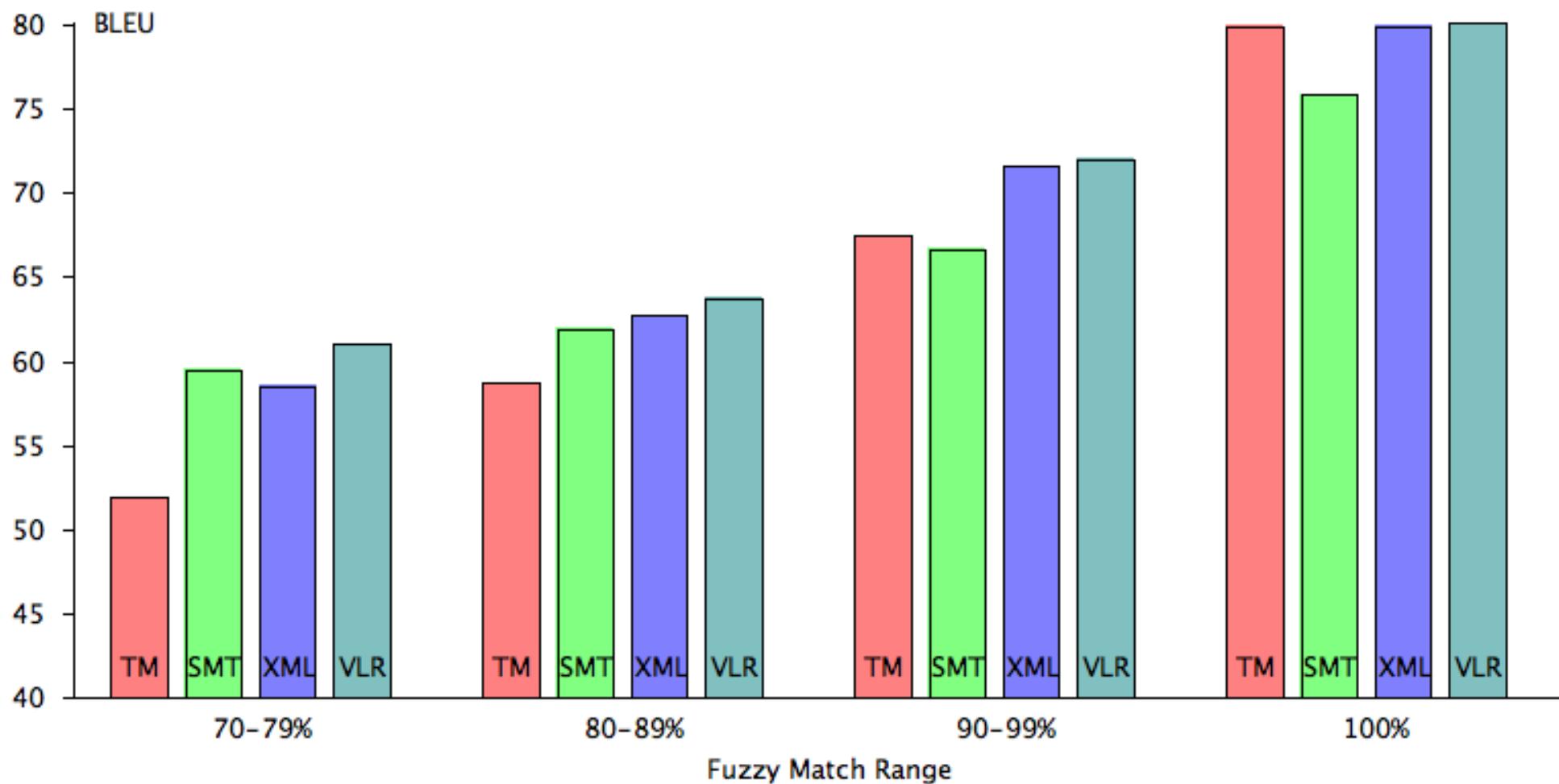
The second paragraph of Article 21 is deleted .

- Very large hierarchical rule

( The second paragraph of Article x is deleted .  
; À l' article x , le texte du deuxième alinéa est supprimé . )

# Result: Acquis

141





# logging and eye tracking





- Different types of events are saved in the logging.
  - configuration and statistics
  - start and stop session
  - segment opened and closed
  - text, key strokes, and mouse events
  - scroll and resize
  - search and replace
  - suggestions loaded and suggestion chosen
  - interactive translation prediction
  - gaze and fixation from eye tracker



- In every event we save:
  - Type
  - In which element was produced
  - Time
- Special attributes are kept for some types of events
  - Diff of a text change
  - Current cursor position
  - Character looked at
  - Clicked UI element
  - Selected text

⇒ Full replay of user session is possible

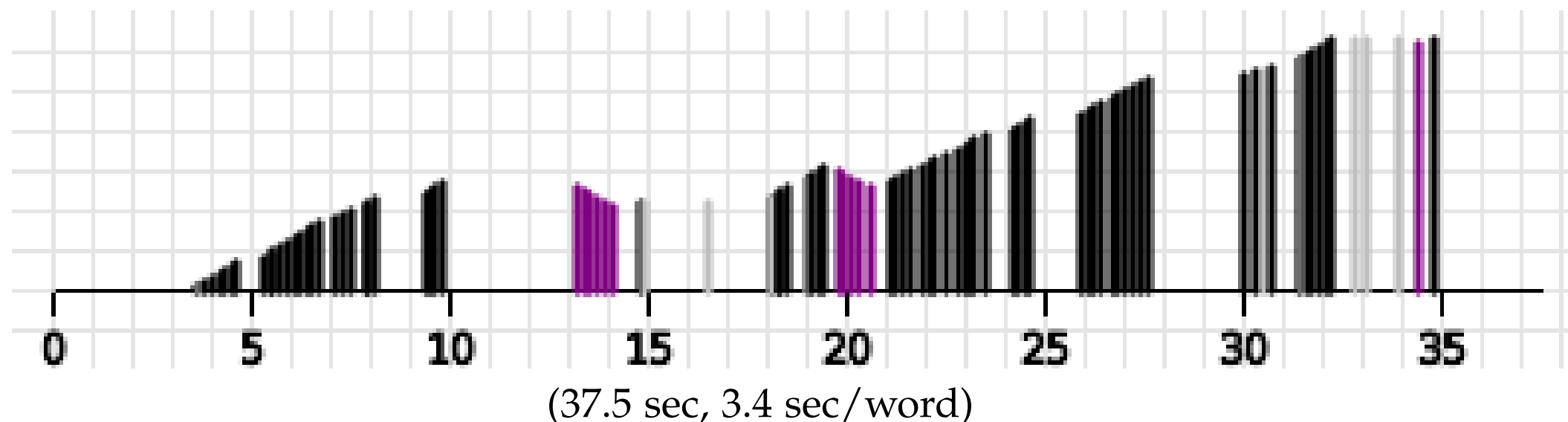
# Keystroke Log

145



Input: Au premier semestre, l'avionneur a livré 97 avions.

Output: The manufacturer has delivered 97 planes during the first half.



black: keystroke, purple: deletion, grey: cursor move  
height: length of sentence

# Example of Quality Judgments

146



---

Src. Sans se démonter, il s'est montré concis et précis.

MT Without dismantle, it has been concise and accurate.

---

1/3 Without fail, he has been concise and accurate. *(Prediction+Options, L2a)*

4/0 Without getting flustered, he showed himself to be concise and precise.  
*(Unassisted, L2b)*

4/0 Without falling apart, he has shown himself to be concise and accurate. *(Postedit, L2c)*

1/3 Unswayable, he has shown himself to be concise and to the point. *(Options, L2d)*

0/4 Without showing off, he showed himself to be concise and precise. *(Prediction, L2e)*

1/3 Without dismantling himself, he presented himself consistent and precise.  
*(Prediction+Options, L1a)*

2/2 He showed himself concise and precise. *(Unassisted, L1b)*

3/1 Nothing daunted, he has been concise and accurate. *(Postedit, L1c)*

3/1 Without losing face, he remained focused and specific. *(Options, L1d)*

3/1 Without becoming flustered, he showed himself concise and precise. *(Prediction, L1e)*

# Main Measure: Productivity

Assistance	Speed	Quality
Unassisted	4.4s/word	47% correct
Postedit	2.7s (-1.7s)	55% (+8%)
Options	3.7s (-0.7s)	51% (+4%)
Prediction	3.2s (-1.2s)	54% (+7%)
Prediction+Options	3.3s (-1.1s)	53% (+6%)

# Faster and Better, Mostly

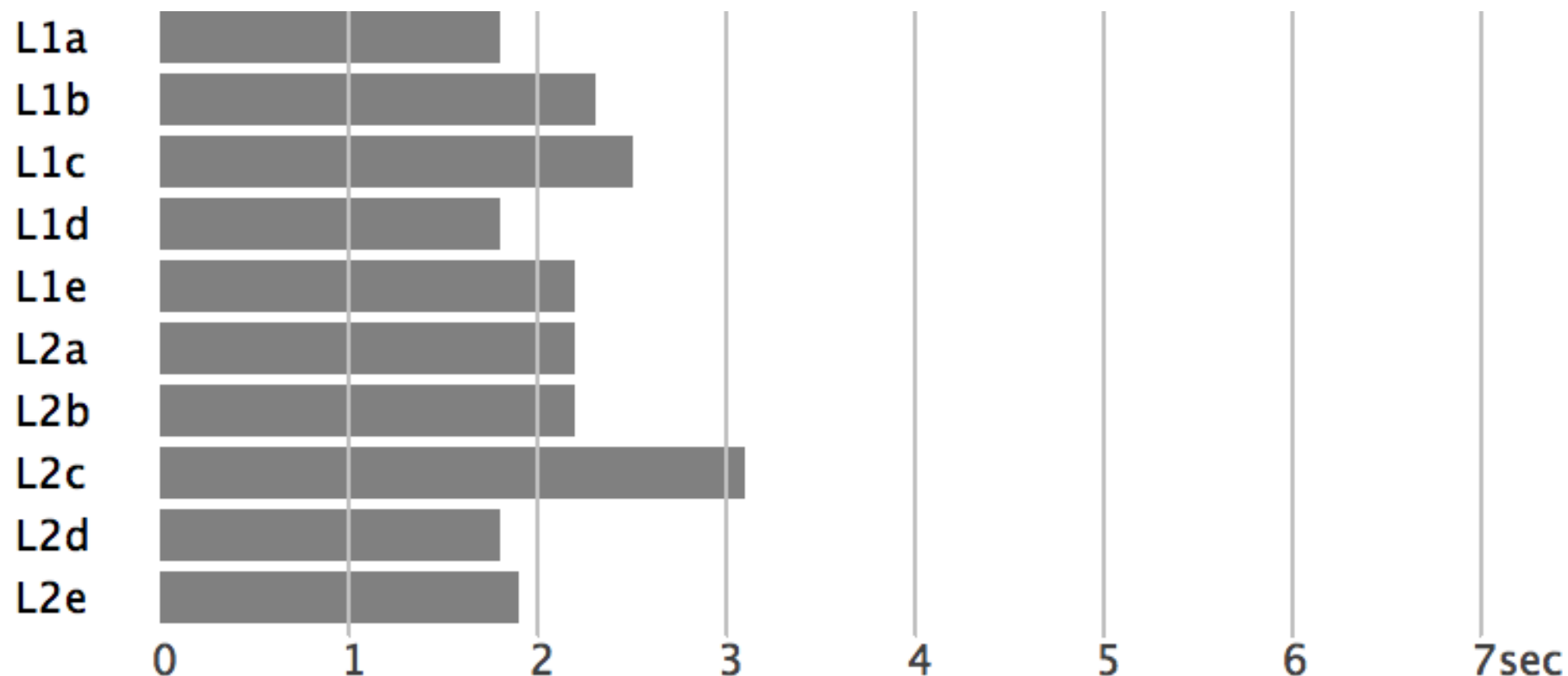
148



User	Unassisted	Postedit		Options		Prediction		Prediction+Options	
L1a	3.3sec/word 23% correct	1.2s 39%	-2.2s +16%)	2.3s 45%	-1.0s +22%	1.1s 30%	-2.2s +7%)	2.4s 44%	-0.9s +21%
L1b	7.7sec/word 35% correct	4.5s 48%	-3.2s +13%	4.5s 55%	-3.3s +20%	2.7s 61%	-5.1s +26%	4.8s 41%	-3.0s +6%
L1c	3.9sec/word 50% correct	1.9s 61%	-2.0s +11%	3.8s 54%	-0.1s +4%	3.1s 64%	-0.8s +14%	2.5s 61%	-1.4s +11%
L1d	2.8sec/word 38% correct	2.0s 46%	-0.7s +8%	2.9s 59%	(+0.1s) (+21%)	2.4s 37%	(-0.4s) (-1%)	1.8s 45%	-1.0s +7%
L1e	5.2sec/word 58% correct	3.9s 64%	-1.3s +6%	4.9s 56%	(-0.2s) (-2%)	3.5s 62%	-1.7s +4%	4.6s 56%	(-0.5s) (-2%)
L2a	5.7sec/word 16% correct	1.8s 50%	-3.9s +34%	2.5s 34%	-3.2s +18%	2.7s 40%	-3.0s +24%	2.8s 50%	-2.9s +34%
L2b	3.2sec/word 64% correct	2.8s 56%	(-0.4s) (-8%)	3.5s 60%	+0.3s -4%	6.0s 61%	+2.8s -3%	4.6s 57%	+1.4s -7%
L2c	5.8sec/word 52% correct	2.9s 53%	-3.0s +1%	4.6s 37%	(-1.2s) (-15%)	4.1s 59%	-1.7s +7%	2.7s 53%	-3.1s +1%
L2d	3.4sec/word 49% correct	3.1s 49%	(-0.3s) (+0%)	4.3s 51%	(+0.9s) (+2%)	3.8s 53%	(+0.4s) (+4%)	3.7s 58%	(+0.3s) (+9%)
L2e	2.8sec/word 68% correct	2.6s 79%	-0.2s +11%	3.5s 59%	+0.7s -9%	2.8s 64%	(-0.0s) (-4%)	3.0s 66%	+0.2s -2%
avg.	4.4sec/word 47% correct	2.7s 55%	-1.7s +8%	3.7s 51%	-0.7s +4%	3.2s 54%	-1.2s +7%	3.3s 53%	-1.1s +6%

# Unassisted Novice Translators

149

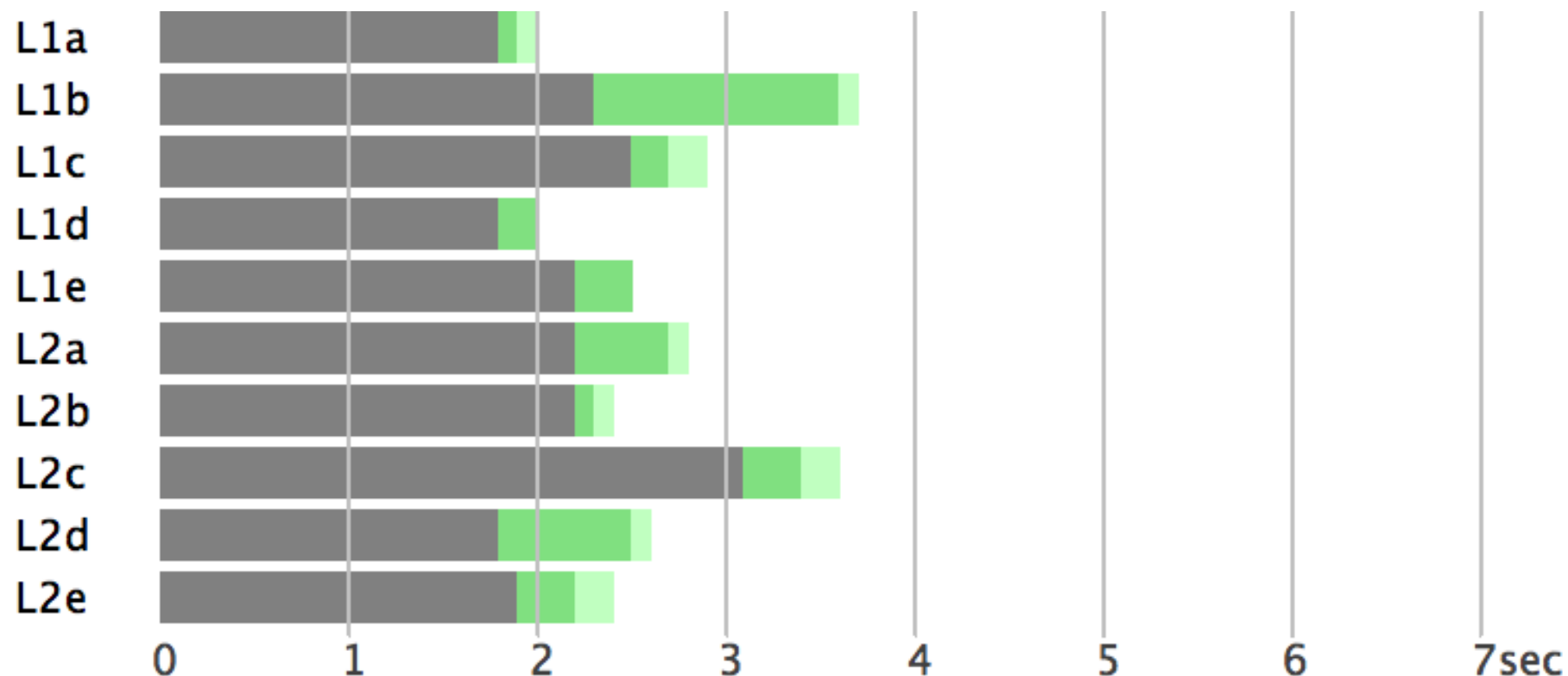


L1 = native French, L2 = native English, average time per input word

only typing

# Unassisted Novice Translators

150



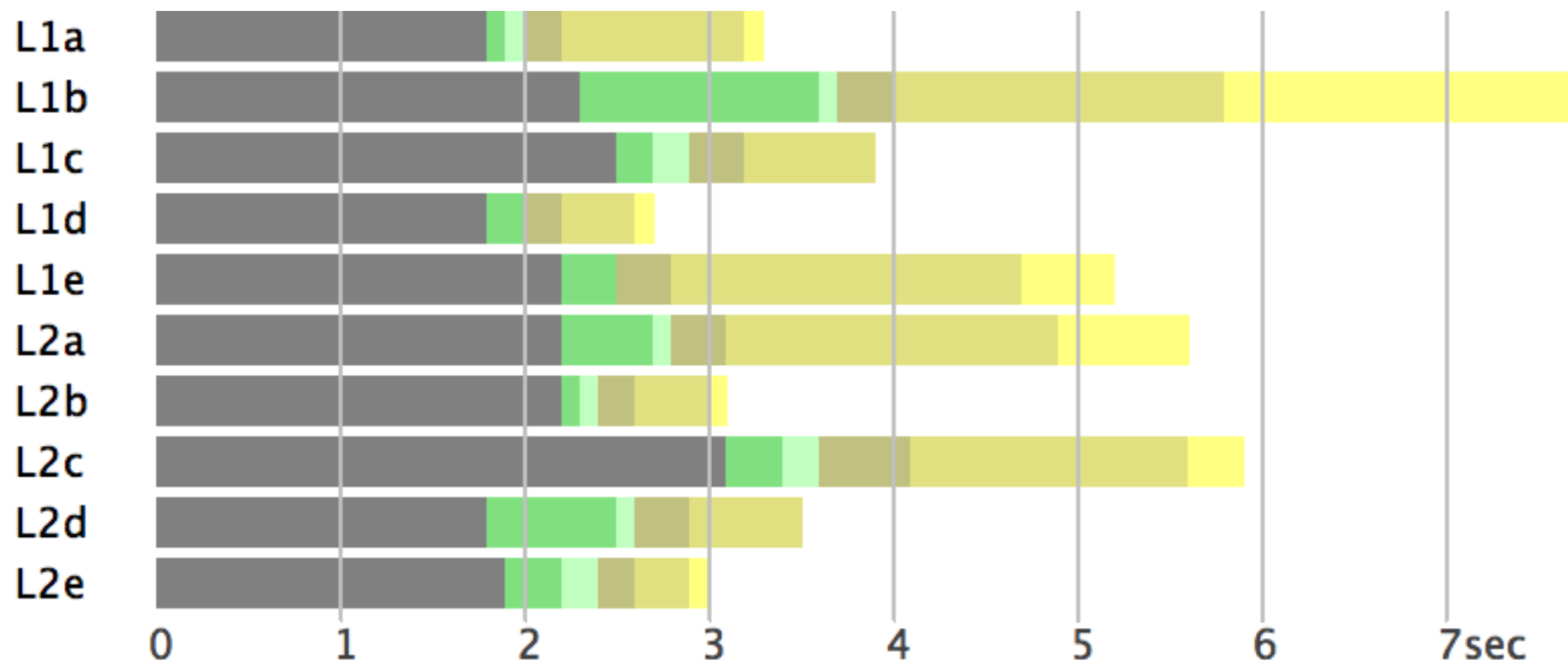
L1 = native French, L2 = native English, average time per input word

typing, initial and final pauses



# Unassisted Novice Translators

151



L1 = native French, L2 = native English, average time per input word

typing, **initial and final pauses**, **short, medium, and long pauses**  
**most time difference on intermediate pauses**

# Activities: Native French User L1b

User: L1b	total	init-p	end-p	short-p	mid-p	big-p	key	click	tab
Unassisted	7.7s	1.3s	0.1s	0.3s	1.8s	1.9s	2.3s	-	-
Postedit	4.5s	1.5s	0.4s	0.1s	1.0s	0.4s	1.1s	-	-
Options	4.5s	0.6s	0.1s	0.4s	0.9s	0.7s	1.5s	0.4s	-
Prediction	2.7s	0.3s	0.3s	0.2s	0.7s	0.1s	0.6s	-	0.4s
Prediction+Options	4.8s	0.6s	0.4s	0.4s	1.3s	0.5s	0.9s	0.5s	0.2s

# Activities: Native French User L1b

153



User: L1b	total	init-p	end-p	short-p	mid-p	big-p	key	click	tab
Unassisted	7.7s	1.3s	0.1s	0.3s	1.8s	1.9s	2.3s	-	-
Postedit	4.5s	1.5s	0.4s	0.1s	1.0s	0.4s	1.1s	-	-
Options	4.5s	0.6s	0.1s	0.4s	0.9s	0.7s	1.5s	0.4s	-
Prediction	2.7s	0.3s	0.3s	0.2s	0.7s	0.1s	0.6s	-	0.4s
Prediction+Options	4.8s	0.6s	0.4s	0.4s	1.3s	0.5s	0.9s	0.5s	0.2s

Slightly  
less time  
spent on  
typing

# Activities: Native French User L1b

154



User: L1b	total	init-p	end-p	short-p	mid-p	big-p	key	click	tab
Unassisted	7.7s	1.3s	0.1s	0.3s	1.8s	1.9s	2.3s	-	-
Postedit	4.5s	1.5s	0.4s	0.1s	1.0s	0.4s	1.1s	-	-
Options	4.5s	0.6s	0.1s	0.4s	0.9s	0.7s	1.5s	0.4s	-
Prediction	2.7s	0.3s	0.3s	0.2s	0.7s	0.1s	0.6s	-	0.4s
Prediction+Options	4.8s	0.6s	0.4s	0.4s	1.3s	0.5s	0.9s	0.5s	0.2s

Less  
pausing

Slightly  
less time  
spent on  
typing

# Activities: Native French User L1b

155



User: L1b	total	init-p	end-p	short-p	mid-p	big-p	key	click	tab
Unassisted	7.7s	1.3s	0.1s	0.3s	1.8s	1.9s	2.3s	-	-
Postedit	4.5s	1.5s	0.4s	0.1s	1.0s	0.4s	1.1s	-	-
Options	4.5s	0.6s	0.1s	0.4s	0.9s	0.7s	1.5s	0.4s	-
Prediction	2.7s	0.3s	0.3s	0.2s	0.7s	0.1s	0.6s	-	0.4s
Prediction+Options	4.8s	0.6s	0.4s	0.4s	1.3s	0.5s	0.9s	0.5s	0.2s

Less  
pausing

Especially  
less time  
in big  
pauses

Slightly  
less time  
spent on  
typing

# Origin of Characters: Native French L1b

156



User: L1b	key	click	tab	mt
Postedit	18%	-	-	81%
Options	59%	40%	-	-
Prediction	14%	-	85%	-
Prediction+Options	21%	44%	33%	-

# Origin of Characters: Native French L1b

157



User: L1b	key	click	tab	mt
Postedit	18%	-	-	81%
Options	59%	40%	-	-
Prediction	14%	-	85%	-
Prediction+Options	21%	44%	33%	-

Translation comes to large degree from assistance



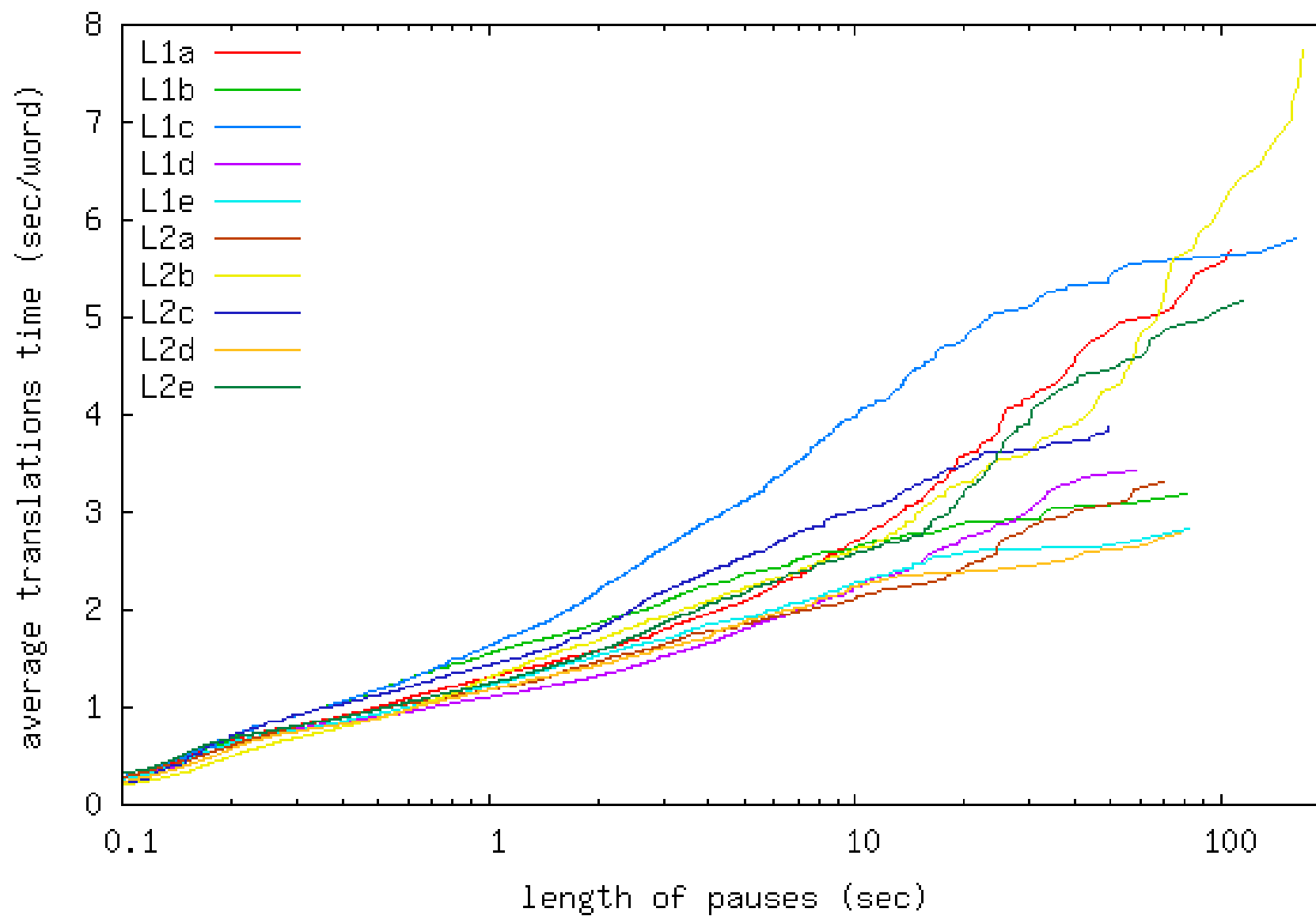
- Our classification of pauses is arbitrary (2-6sec, 6-60sec, >60sec)
- Extreme view: all you see is pauses
  - keystrokes take no observable time
  - all you see is pauses between action points■
- Visualizing range of pauses:  
time  $t$  spent in pauses  $p \in P$  up to a certain length  $l$

$$sum(t) = \frac{1}{Z} \sum_{p \in P, l(p) \leq t} l(p)$$



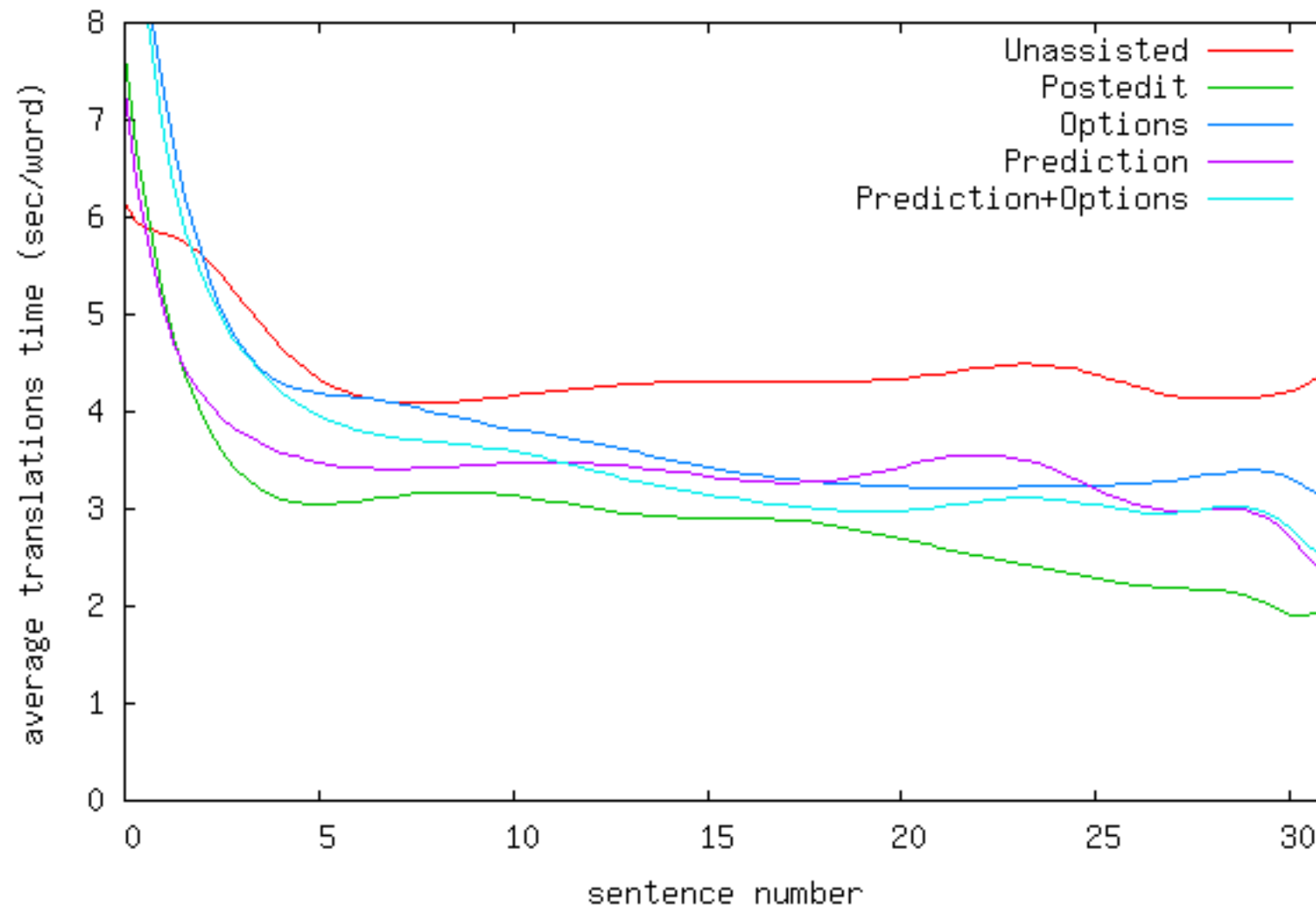
# Results

159

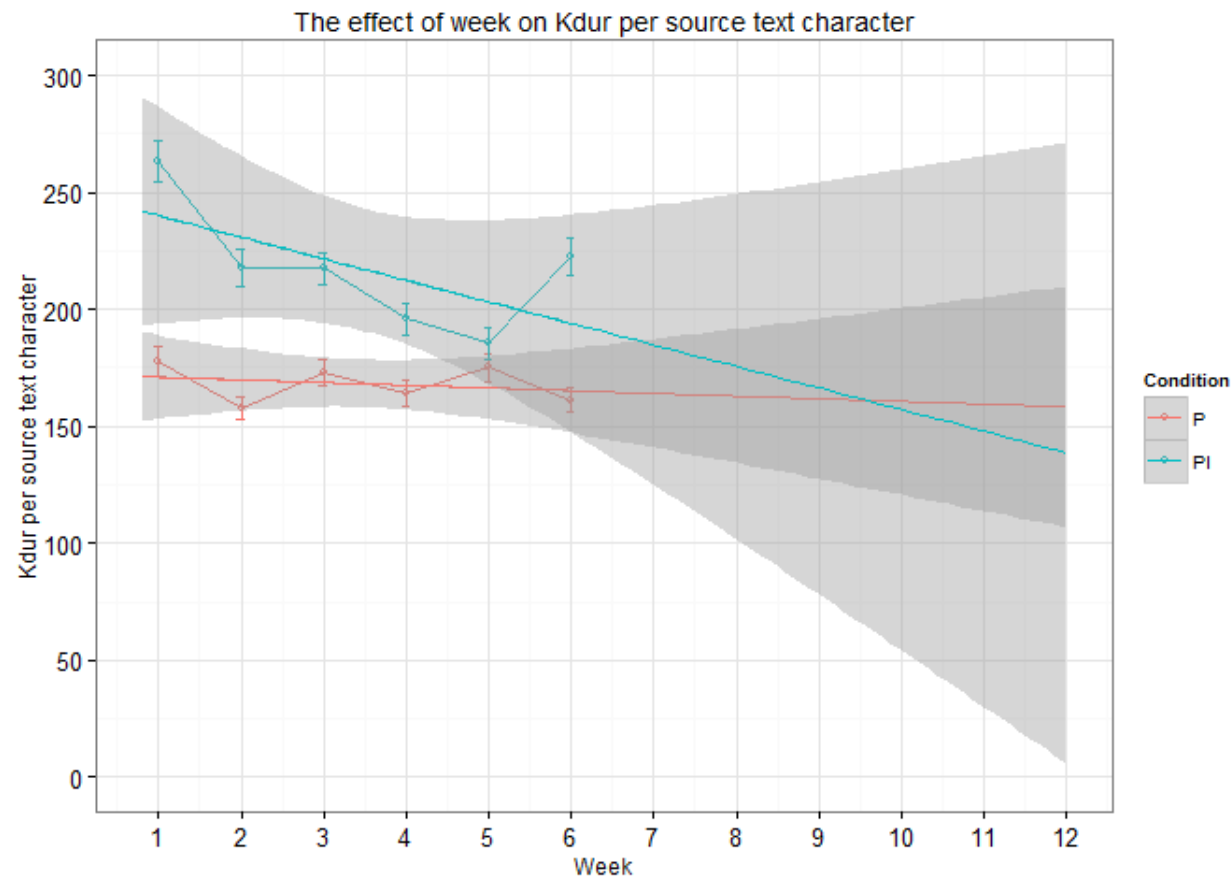


# Learning Effects

Users become better over time with assistance



# Learning Effects: Professional Translators 161



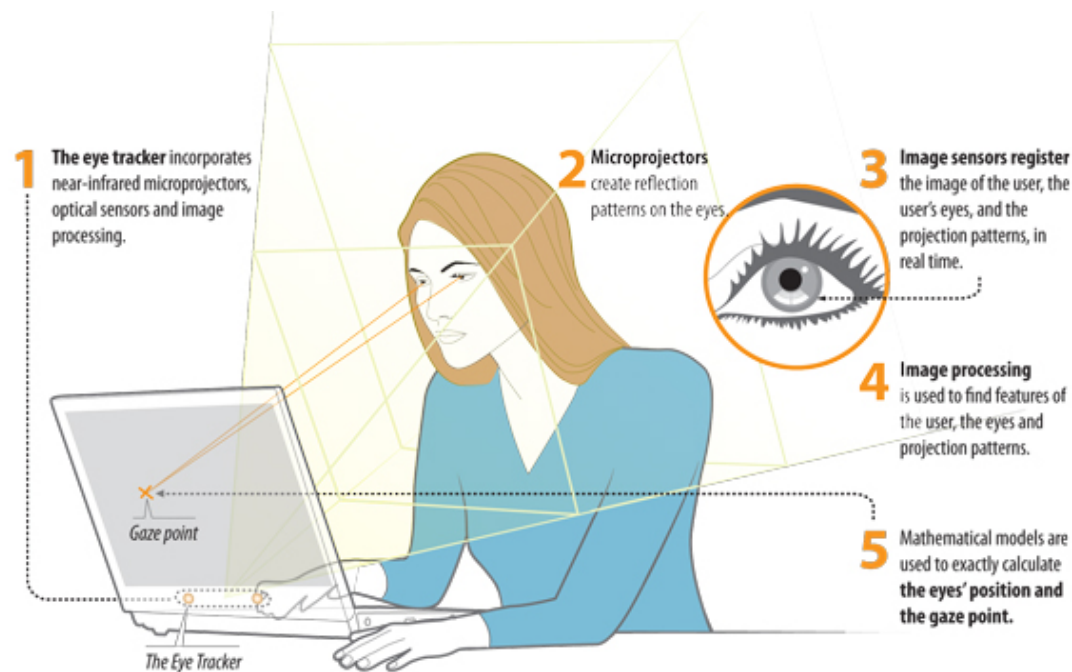
CASMACAT longitudinal study

Productivity projection as reflected in Kdur taking into account six weeks

(Kdur = user activity excluding pauses > 5 seconds)

# Eye Tracking

162




- Eye trackers extensively used in cognitive studies of, e.g., reading behavior
- Overcomes weakness of key logger: what happens during pauses
- Fixation: where is the focus of the gaze
- Pupil dilation: indicates degree of concentration



- Problem: Accuracy and precision of gaze samples




Good precision,  
poor accuracy

The diagram shows a black bullseye target. A cluster of approximately 10 red 'x' marks, representing eye tracker results, is tightly grouped together in the upper right quadrant of the target, indicating high precision but low accuracy.

Good accuracy,  
poor precision

The diagram shows a black bullseye target. Red 'x' marks are scattered across the entire area of the target, with a few marks near the center, indicating low precision but high accuracy.

*x = eye tracker result*

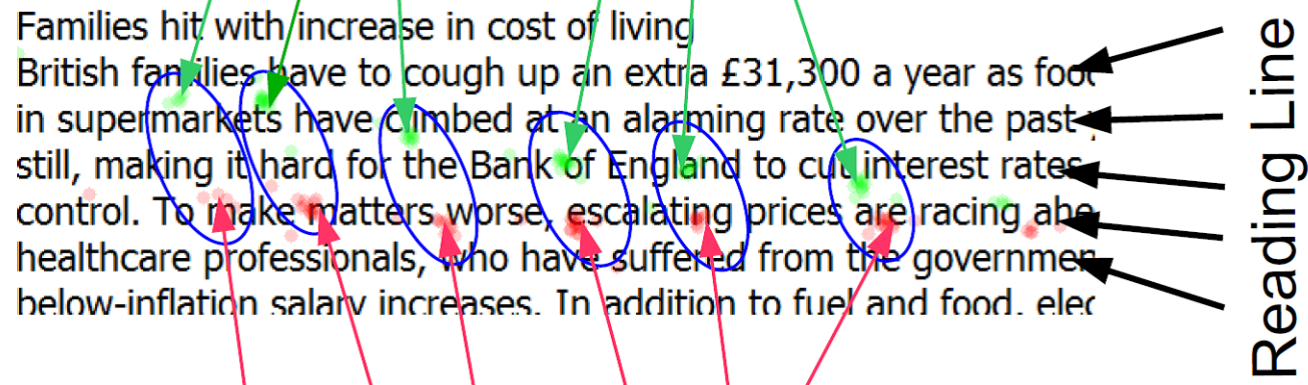
 *= target looked at*

# Gaze-to-Word Mapping

- Recorded gaze locations and fixations

## Right eye gaze samples

Families hit with increase in cost of living  
British families have to cough up an extra £31,300 a year as food  
in supermarkets have climbed at an alarming rate over the past  
still, making it hard for the Bank of England to cut interest rates  
control. To make matters worse, escalating prices are racing ahead  
healthcare professionals, who have suffered from the government  
below-inflation salary increases. In addition to fuel and food, elec

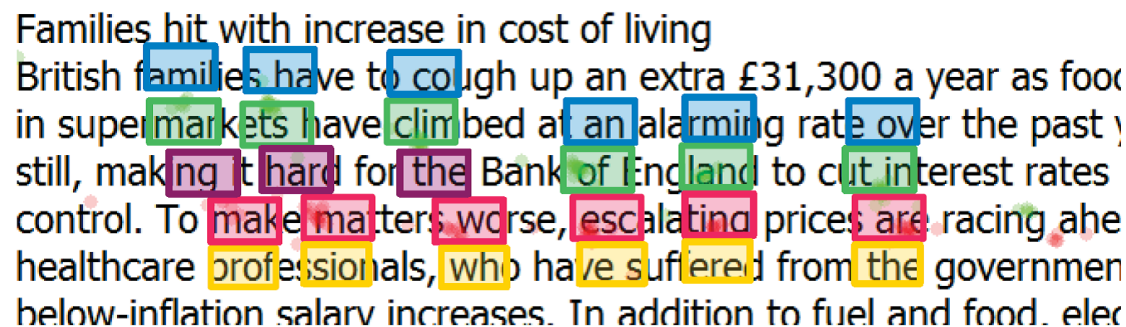


Reading Line

## Left eye gaze samples

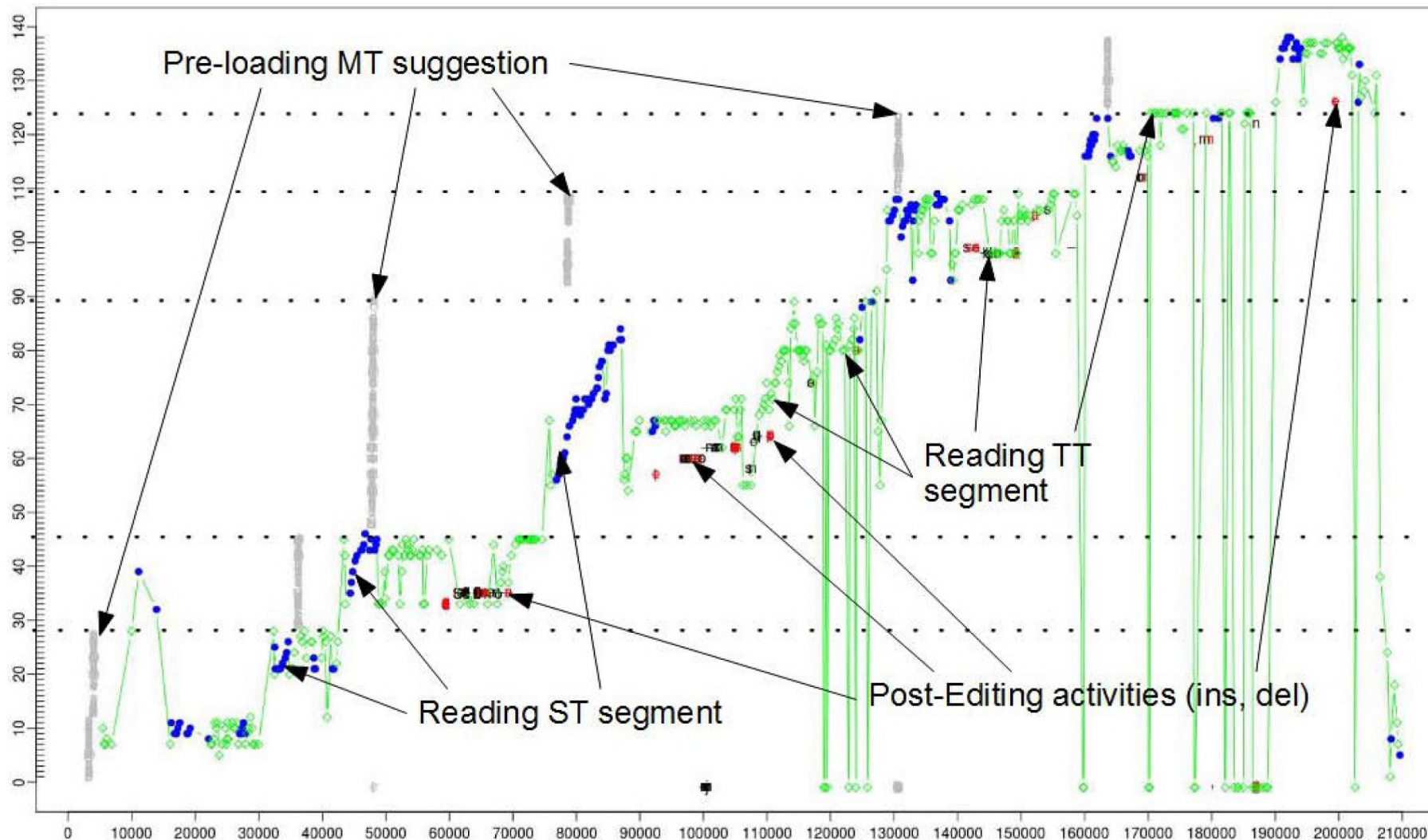
- Gaze-to-word mapping

Families hit with increase in cost of living  
British families have to cough up an extra £31,300 a year as food  
in supermarkets have climbed at an alarming rate over the past  
still, making it hard for the Bank of England to cut interest rates  
control. To make matters worse, escalating prices are racing ahead  
healthcare professionals, who have suffered from the government  
below-inflation salary increases. In addition to fuel and food, elec



# Logging and Eye Tracking

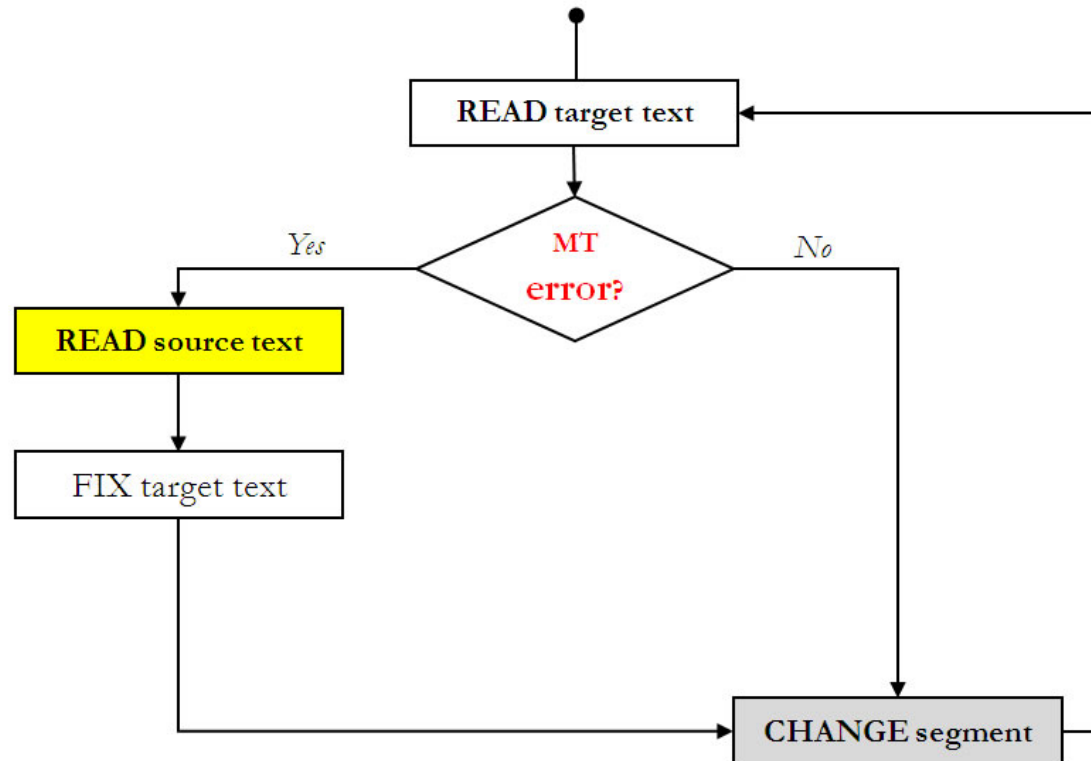
165



focus on target word (green) or source word (blue) at position  $x$

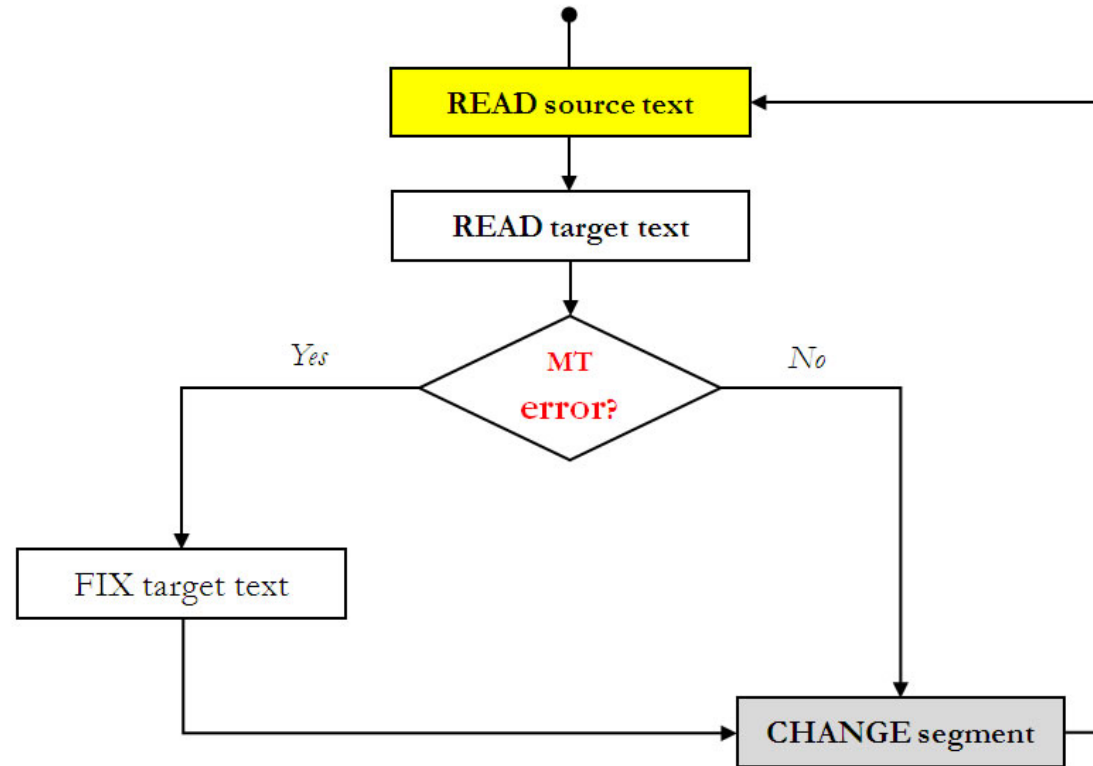
# Cognitive Studies: User Styles

166

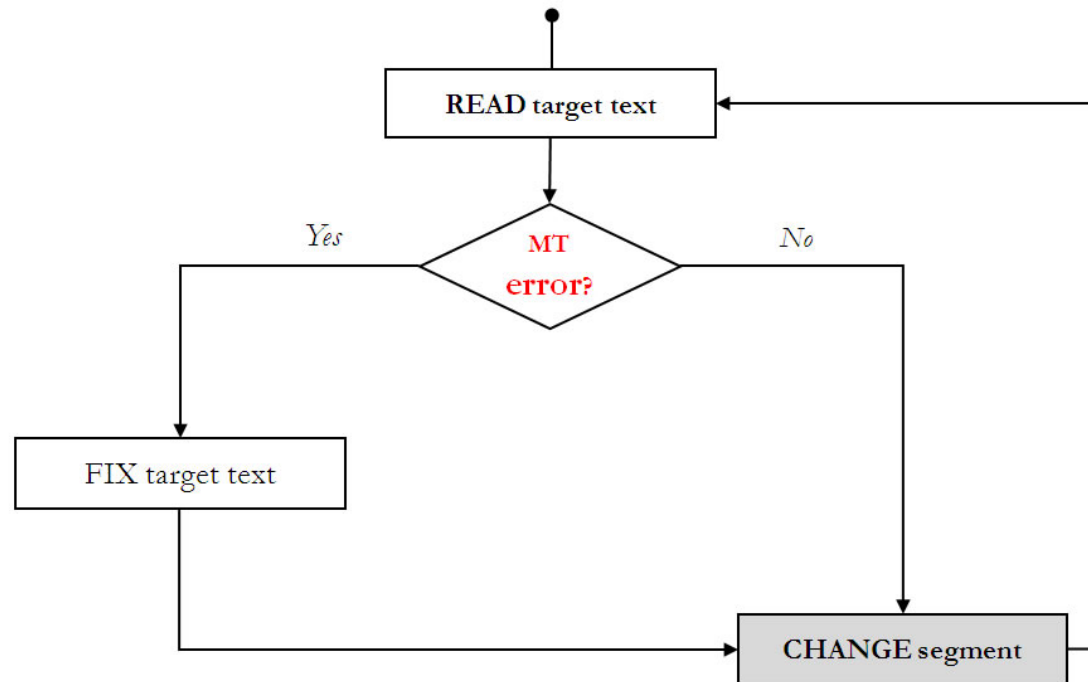


- User style 1: Verifies translation just based on the target text, reads source text to fix it

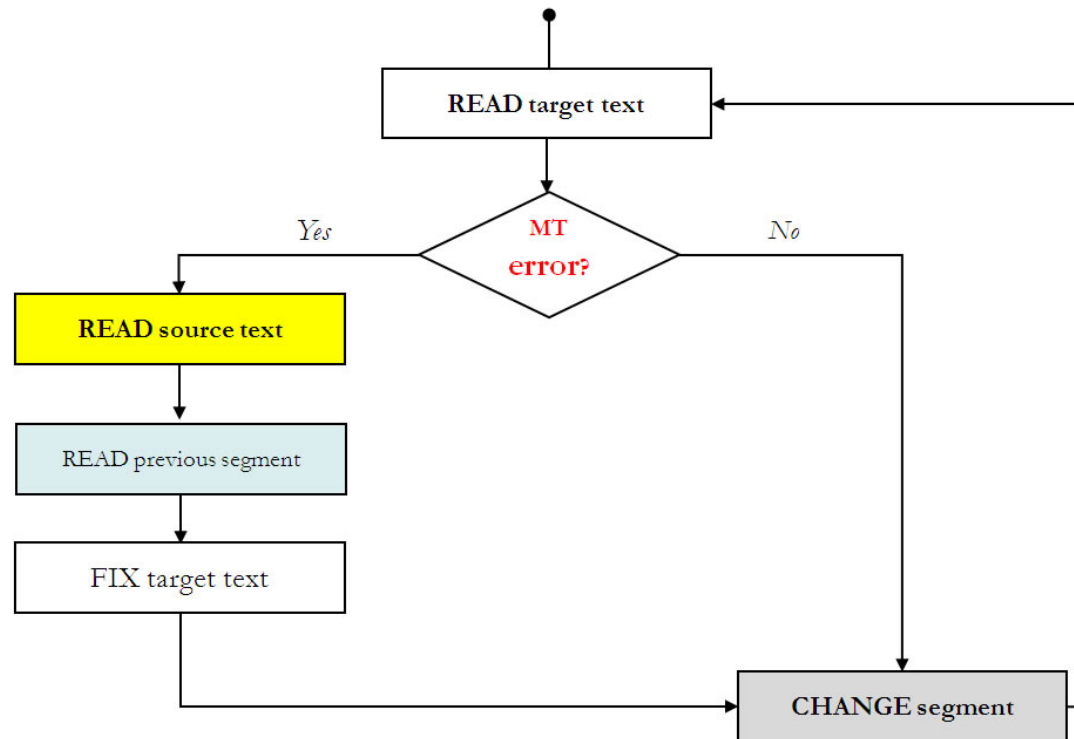




- User style 2: Reads source text first, then target text



- User style 3: Makes corrections based on target text only



- User style 4: As style 1, but also considers previous segment for corrections

# Users and User Styles

	Style 1			Style 2			Style 3			Style 4		
	target / source-fix			source-target			target only			wider context		
	P	PI	PIA	P	PI	PIA	P	PI	PIA	P	PI	PIA
P02	*	*	*	●	●	●	●			●	●	●
P03												
P04	●	*	*				*	●	●	●	●	●
P05	●	●	●				*	*	*	●	●	●
P07	*	*	*				●	●	●	●	●	●
P08	*	*	*	●	●	●				●	●	●
P09	●	●	●				*	*	*	●	●	●

- Individual users employ different user styles
- But: consistently across different types of assistance  
(P = post-editing, PI = interactive post-editing, PIA = interactive post-editing with additional annotations)



- Local backtracking
  - **Immediate repetition:** the user immediately returns to the same segment (e.g. AAAA)
  - **Local alternation:** user switches between adjacent segments, often singly (e.g. ABAB) but also for longer stretches (e.g. ABC-ABC).
  - **Local orientation:** very brief reading of a number of segments, then returning to each one and editing them (e.g. ABCDE-ABCDE).
- Long-distance backtracking
  - **Long-distance alternation:** user switches between the current segment and different previous segments (e.g. JCJDJFJG)
  - **Text final backtracking:** user backtracks to specific segments after having edited all the segments at least once
  - **In-text long distance backtracking:** instances of long distance backtracking as the user proceeds in order through the text.

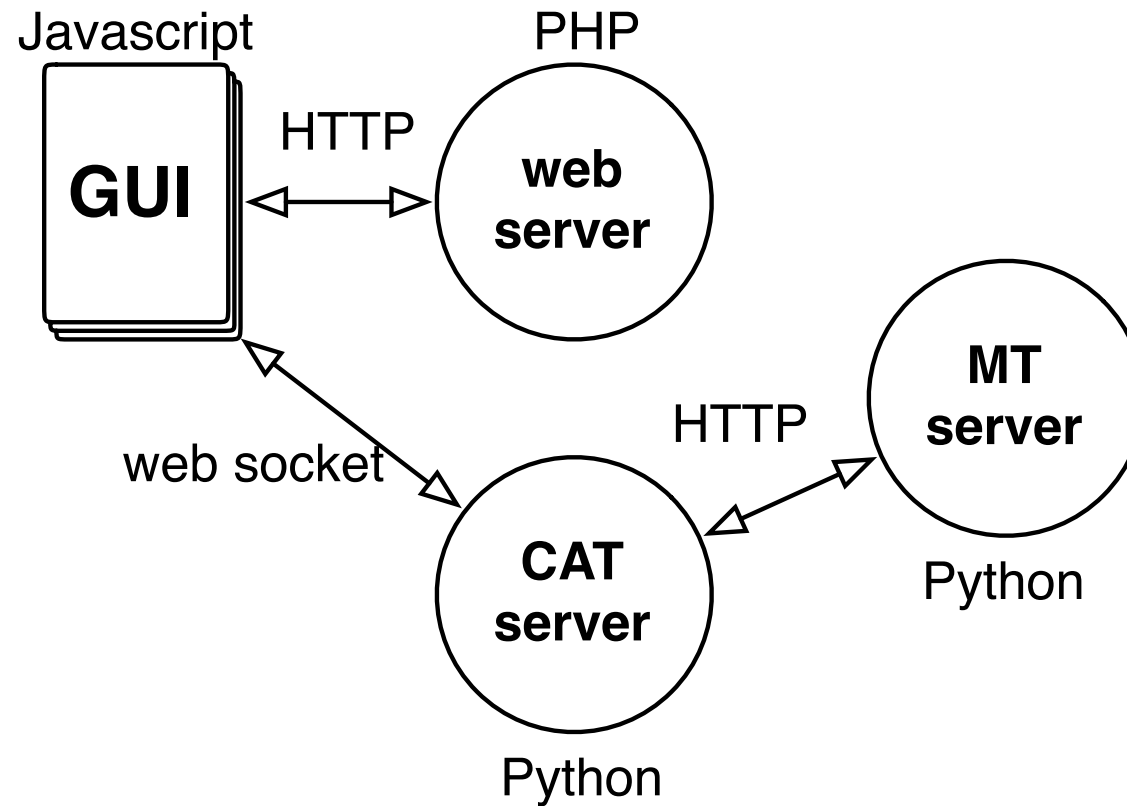


# part III

## CASMACAT workbench implementation

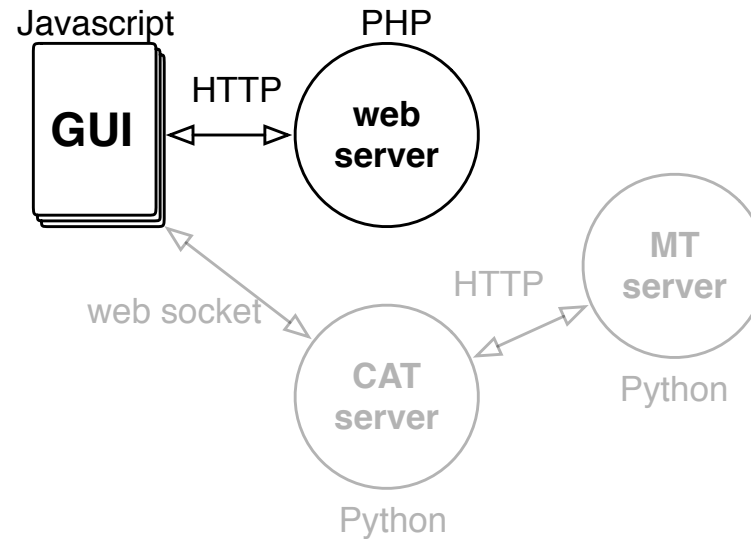
# Components

173



# Web Server

174



- Builds on Matecat open source implementation
- Typical web application: LAMP (Linux, Apache, MySQL, PHP)
- Uses model, view, controller breakdown

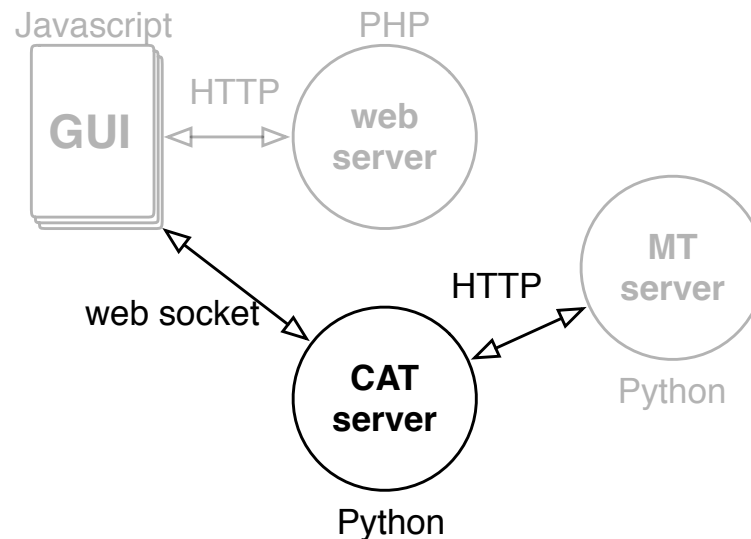




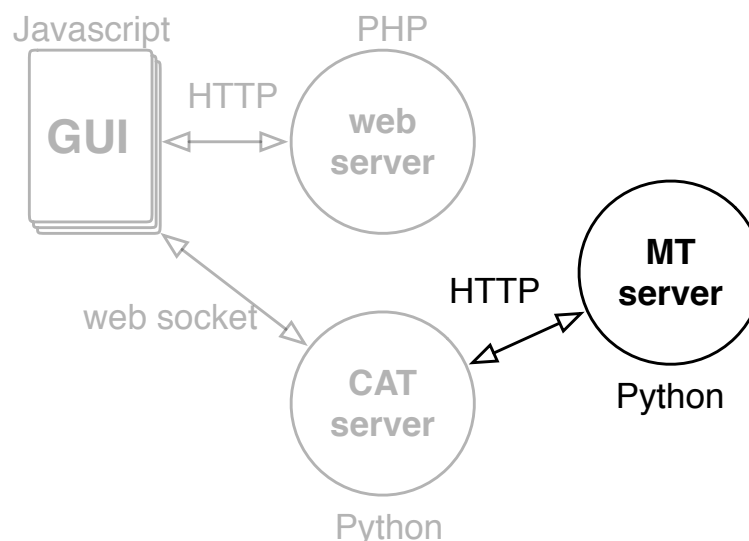
- Relevant data is stored in MySQL database `matecat_sandbox`
- Major database tables
  - Projects are stored in `projects`
  - They have a corresponding entry in `jobs`
  - Raw files (XLIFF) are stored in `files`
  - Segments are stored in `segments`
  - Translations of segments are stored in `segment_translations`
  - Log events are stored in `*_event`
  - etc.
- The major change from Matecat is the logging



- Typical request: get information about a segment:  
POST `http://192.168.56.2:8000/?action=getSegments&time=1446185242727`
- Script `index.php` selects corresponding action in `lib/controller`  
e.g., `getSegmentsController.php`
- Response is HTML or JSON
- The main action is really in the Javascript GUI `public/js`
  - core functionality from Matecat `public/js/cat.js`
  - CASMACAT extensions `public/js/casmacat`



- To a large degree middleware
- Calls external services such as
  - MT server
  - word aligner
  - interactive translation prediction
- Caches information about a sentence translation



- Google-style API to MT Server
- Python wrapper for Moses
  - basic translation request
  - includes pre and post processing pipeline
  - other functions: word alignment, incremental updating, etc.
- Uses mosesserver XMLRPC server

- Requires mosesserver to run as a service

```
mosesserver -config $MODELDIR/moses.ini --server-port 9010
```

- Script server.py requires a lot of parameters
  - preprocessing tools (tokenizer, truecaser, etc.)
  - IP address and port
  - URL of the mosesserver API
  - etc.

- Request to the script

```
http://127.0.0.1:9000//translate?q=Un+test&key=0&source=xx&target=xx
```

- Response

```
{"data": {"translations": [{"translatedText": "A test",  
"translatedTextRaw": "a test",  
"annotatedSource": "un test",  
"tokenization": {"src": [[0, 1], [3, 6]], "tgt": [[0, 0], [2, 5]]}}]}
```



- Moses is installed in `/opt/moses`
- CASMACAT is installed in `/opt/casmacat`
  - web server / GUI in `/opt/casmacat/web-server`
  - MT server (server.py) in `/opt/casmacat/mt-server`
  - CAT server in `/opt/casmacat/cat-server`
  - installation scripts in `/opt/casmacat/install`
  - log files in `/opt/casmacat/logs`
- Home Edition
  - admin web server in `/opt/casmacat/admin`
  - corpus data in `/opt/casmacat/data`
  - prototype training in `/opt/casmacat/experiment`
  - engines stored in `/opt/casmacat/engines`

# Home Edition MT Engine

- Demo engine in /opt/casmacat/engines/fr-en-upload-1
- Files
  - biconcor.1
  - biconcor.1.align
  - biconcor.1.src-vcb
  - biconcor.1.tgt
  - biconcor.1.tgt-vcb
  - corpus-1.binlm.1
  - fast-align.1
  - fast-align.1.log
  - fast-align.1.parameters
  - fast-align-inverse.1
  - fast-align-inverse.1.log
  - fast-align-inverse.1.parameters
  - info
  - moses.tuned.ini.1
  - phrase-table-mmsapt.1
  - reordering-table.1.wbe-msd-bidirectional-fe.minlexr
  - RUN
  - truecase-model.1.en
  - truecase-model.1.fr
- The script RUN starts the engine

# Thank You

182



# questions?